

Smart Intelligent Energy Insight System (SIEIS)

Course: AAI-530

Project Team: Aishwarya Gulhane (Project Lead), Bhavleen Kaur, Paul Ittoopunny

Instructor: Haisav Chokshi

Date: February 26, 2026

Abstract

The Smart Intelligent Energy Insight System (SIEIS) is a distributed data engineering and analytics platform for streaming environmental telemetry. The system integrates event-driven ingestion, dual-path persistence, API-based serving, dashboard-based observability, and scheduled model lifecycle operations. The objective is to support both low-latency operational monitoring and long-term analytical and ML workflows using a unified architecture.

1. Introduction

Indoor environments influence comfort, productivity, and health. Modern IoT deployments stream temperature, humidity, ambient light, and device voltage from distributed nodes. Practical systems must support heterogeneous data rates, near real-time visibility, and durable historical retention for analytics and model training. SIEIS addresses these needs by combining stream processing, time-series storage, object storage, and a lightweight ML pipeline in a containerized architecture.

1.1 Project Objectives

- Implement an end-to-end streaming pipeline for environmental telemetry with reliable ingestion and transformation.
- Support a dual persistence strategy: hot storage for low-latency queries and cold storage for durable analytics.
- Expose APIs for health checks, time-series queries, and machine learning inference.
- Provide a dashboard for operational visibility and exploratory analytics.
- Automate simulator management and periodic model retraining.

2. System Architecture

2.1 Architectural Pattern

SIEIS follows a dual-write architecture. The hot path persists readings into InfluxDB to enable low-latency monitoring and interactive queries. The cold path archives readings into MinIO as Parquet files to support durable storage, batch analytics, and model training. A serving layer (FastAPI) and observability layer (Streamlit) consume these stores, while a scheduler coordinates simulator restarts and model lifecycle operations.

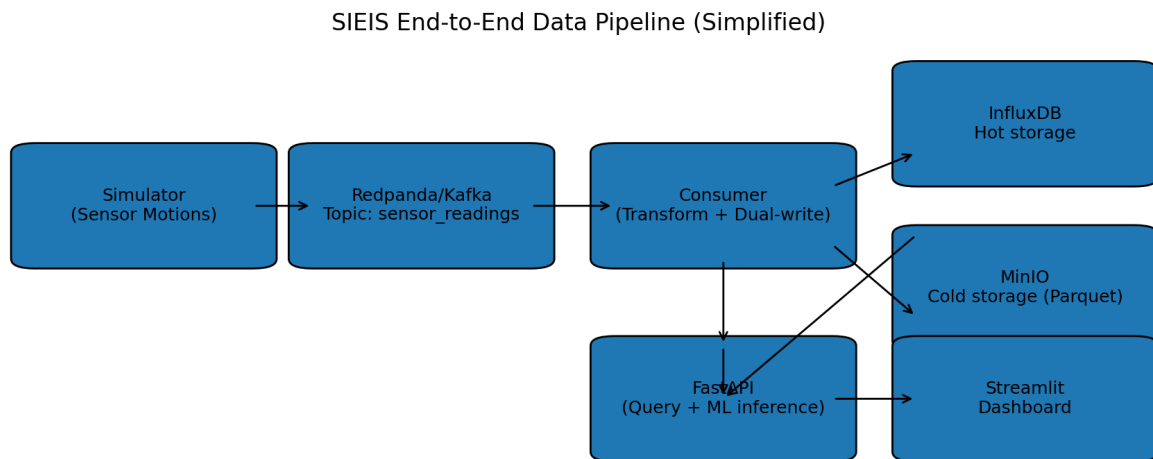


Figure 1. Simplified end-to-end SIEIS data pipeline.

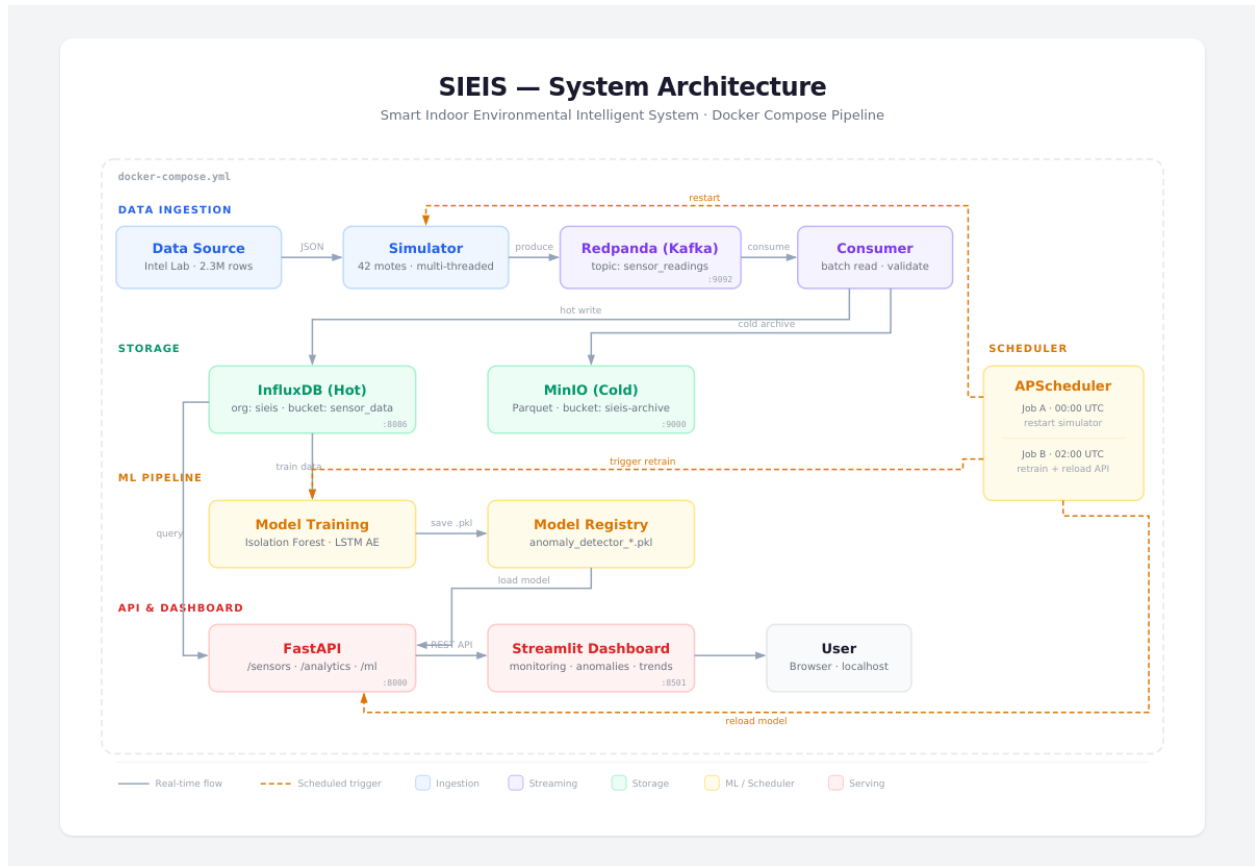


Figure 2. SIEIS system architecture (from project repository).

2.2 Services and Exposed Ports

Table 1 lists the containerized services and their externally exposed ports when running via Docker Compose.

Service	Container	Host Port(s)	Internal Address	Verification
Redpanda (Kafka broker)	sieis-redpanda	19092 (host clients), 9092	redpanda:9092	docker exec sieis-redpanda rpk cluster info
Redpanda Console	sieis-console	8080	redpanda-console:8080	Open http://localhost:8080
InfluxDB	sieis-influxdb3	8086	influxdb3:8086	GET http://localhost:8086/health
MinIO API	sieis-minio	9000	minio:9000	GET http://localhost:9000/minio/health/live
MinIO Console	sieis-minio	9001	minio:9001	Open http://localhost:9001
FastAPI	sieis-api	8000	sieis-api:8000	GET http://localhost:8000/api/v1/health

				alth
Streamlit Dashboard	sieis-dashboard	8501	sieis-dashboard:8501	Open http://localhost:8501
Scheduler	sieis-scheduler	—	internal only	docker logs sieis-scheduler

Table 1. SIEIS services and ports (Docker Compose).

3. Data Description and Preparation

SIEIS uses historical environmental telemetry (originally timestamped in 2004) transformed into a real-time compatible timeline for ingestion and hot-storage acceptance. Each record contains: original date and time, epoch identifier, mote ID, and sensor readings (temperature, humidity, light, voltage). A derived `updated_timestamp` maps original timestamps into a recent time window to prevent future-dated writes.

3.1 Dataset Summary

Records (historical_data)	1,850,934
Unique motes	58
Updated timestamp range (UTC)	2026-01-31 19:03 to 2026-02-23 11:30
Primary sensor features	temperature, humidity, light, voltage
ML time features	hour, day_of_week

Table 2. Dataset summary used for analysis in this report.

3.2 Exploratory Data Analysis

Exploratory analysis was performed on a random sample after applying basic sensor bounds consistent with physical plausibility. Figures 3–5 summarize typical temporal behavior, marginal distributions, and pairwise correlations.

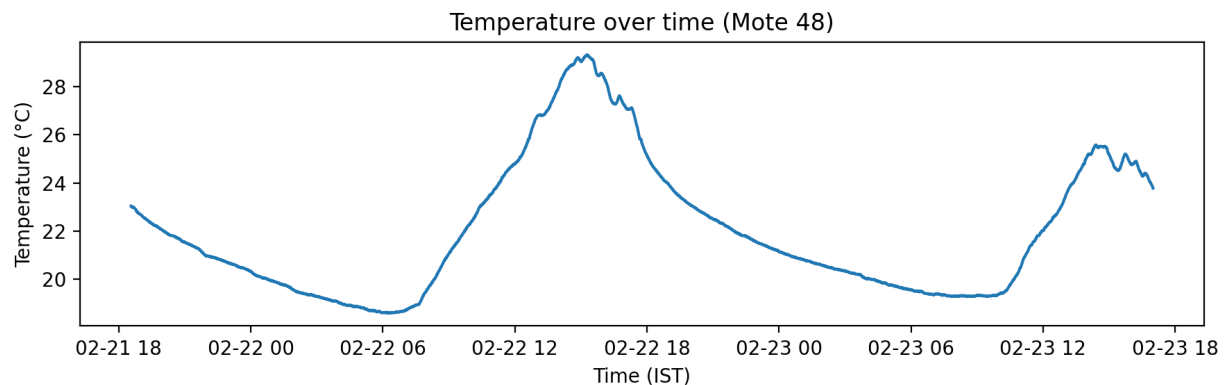


Figure 3. Temperature time series example for a single mote (IST).

Sensor Value Distributions (Sample)

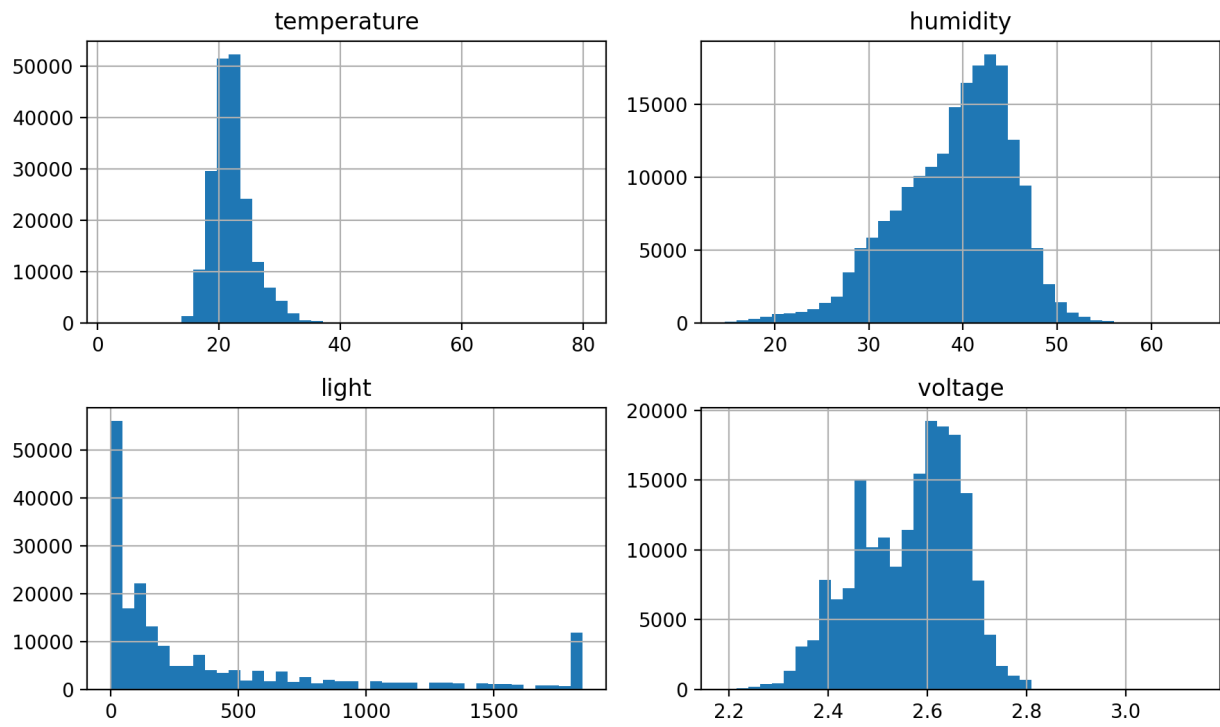


Figure 4. Marginal distributions of sensor features (sample).

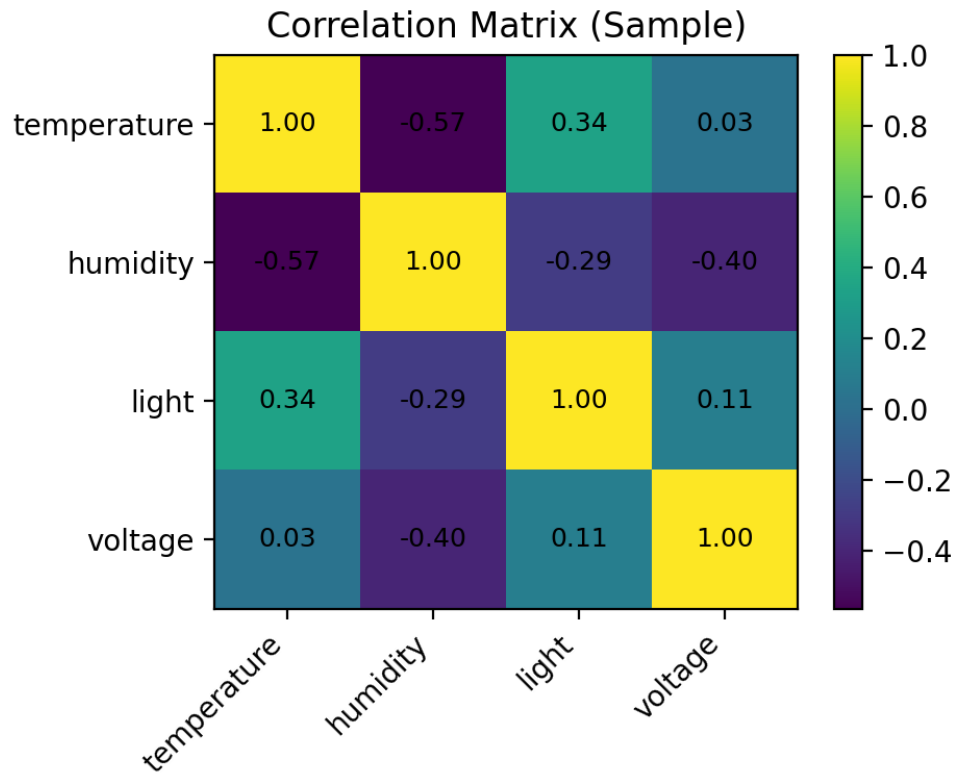


Figure 5. Correlation matrix for sensor features (sample).

4. Machine Learning: Anomaly Detection

The system applies unsupervised anomaly detection to identify unusual readings that may indicate sensor faults, environmental events, or data quality issues. With no anomaly labels available, an Isolation Forest model learns the normal operating region of the feature space and assigns anomaly scores to new observations.

4.1 Feature Engineering

The anomaly detector uses sensor features (temperature, humidity, light, voltage) and time-derived features (hour of day and day of week) computed from `updated_timestamp`. Missing values are imputed using feature medians after range filtering.

4.2 Training Configuration

An Isolation Forest is trained with `StandardScaler` preprocessing. The `contamination` parameter is set to 0.05 (expected 5% outlier fraction). Training was performed on a 200,000-row sample for computational efficiency.

Training samples	188,013
Features	temperature, humidity, light, voltage, hour, day_of_week
Contamination	0.05
Trees (n_estimators)	150
Anomalies detected	9,401 (5.00%)
Decision score mean \pm std	0.0789 \pm 0.0404

Table 3. Isolation Forest training configuration and summary metrics.

4.3 Results and Visualization

Figure 6 shows the distribution of decision scores (higher indicates more normal behavior). Figure 7 overlays detected anomalies on a mote’s temperature series to illustrate how the model highlights unusual points in context.

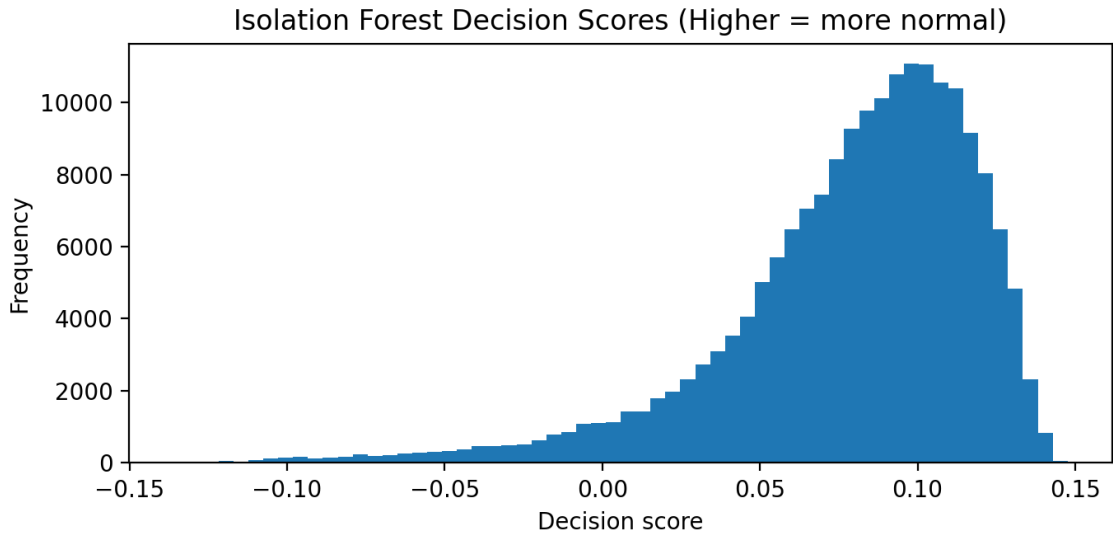


Figure 6. Distribution of Isolation Forest decision scores (sample).

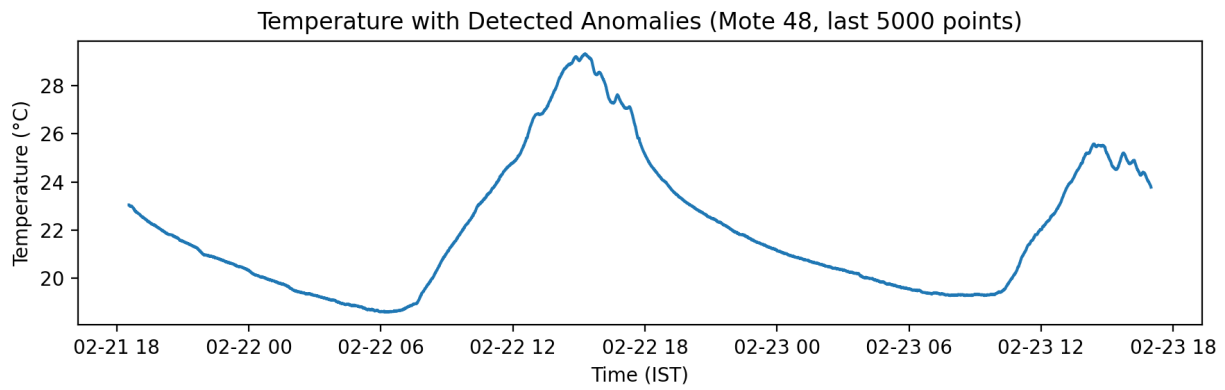


Figure 7. Example anomalies highlighted on temperature (single mote).

5. Deployment, Observability, and Operations

SIEIS is deployed locally via Docker Compose, including the simulator, message broker, consumer, hot and cold stores, API server, dashboard, and scheduler. Operational health is verified via service health endpoints, broker tooling, and container logs. The Streamlit dashboard provides live telemetry visualization and supports interactive inspection of stored data and ML inference results.

5.1 Automation and Model Lifecycle

A scheduler service coordinates recurring tasks such as restarting the simulator and retraining the anomaly detector. After training, the model artifact is saved and can be reloaded into the API service for online inference, enabling a practical MLOps loop within the same containerized environment.

6. Limitations and Future Work

- Ground-truth anomaly labels are not available; evaluation relies on qualitative inspection and score distributions.
- The current model treats each record independently; future work could incorporate temporal models (e.g., autoencoders) for sequence anomalies.
- Per-mote personalization could improve sensitivity by training separate models or adding mote identifiers as features.
- Security hardening is required for production: secret management, authentication, TLS, and role-based access control.

7. Conclusion

SIEIS demonstrates an end-to-end IoT data engineering and analytics platform integrating streaming ingestion, dual-path storage, API-based serving, dashboard observability, and scheduled ML lifecycle operations. The architecture supports both real-time monitoring and durable analytics, while the anomaly detection pipeline provides a practical mechanism to surface unusual sensor behavior without labeled data.

References

1. Liu, F. T., Ting, K. M., & Zhou, Z.-H. (2008). Isolation Forest. In Proceedings of the 2008 IEEE International Conference on Data Mining.
2. Kreps, J., Narkhede, N., & Rao, J. (2011). Kafka: A Distributed Messaging System for Log Processing. Proceedings of NetDB.
3. InfluxData. (n.d.). InfluxDB documentation. Retrieved 2026-02-26.
4. MinIO, Inc. (n.d.). MinIO documentation. Retrieved 2026-02-26.

Appendix A: Code Base

This appendix contains the consolidated code base used to implement the Smart Intelligent Energy Insight System (SIEIS). Listings are organized by repository-relative path.

File	Type	Size (KB)
.env.example	example	2.1
Dockerfile.api	api	0.2
Dockerfile.consumer	consumer	0.4
Dockerfile.dashboard	dashboard	0.3
Dockerfile.scheduler	scheduler	1.6
Dockerfile.simulator	simulator	0.4
Documentation/QUICKSTART.md	md	2.9
README.md	md	14.2
data/.checkpoint_historical.json	json	0.2
data/raw/mote_locs.txt	txt	0.5
data/realtime_mapping/README.md	md	6.7
data/realtime_mapping/preview_mapping.py	py	5.3
data/realtime_mapping/transform_to_realtime.py	py	8.9
data/realtime_mapping/validate_output.py	py	6.7
docker-compose.yml	yml	8.5
pytest.ini	ini	0.4
requirements.txt	txt	0.5
scripts/check_kafka_messages.py	py	0.7
scripts/check_latest_timestamps.py	py	1.2
scripts/clear_storage.py	py	7.0
scripts/compact_minio_parquet.py	py	11.1
scripts/load_historical_data.py	py	19.0
scripts/load_historical_data_new.py	py	10.9
scripts/ml/diagnose_model.py	py	2.8
scripts/ml/diagnose_registry.py	py	0.2
scripts/remap_timestamps.py	py	11.1
scripts/retrain_model.py	py	4.2
scripts/run_dashboard.py	py	0.7
scripts/simple_verify.py	py	2.2
scripts/split_dataset.py	py	7.6
scripts/train_model.py	py	7.9
scripts/validate_data_quality.py	py	5.5
scripts/verify_influxDb.py	py	5.4
scripts/verify_minio_storage.py	py	2.1
src/_init_.py	py	0.0
src/app/_init_.py	py	0.0
src/app/api/_init_.py	py	0.0
src/app/api/main.py	py	2.4
src/app/api/routes/_init_.py	py	0.0
src/app/api/routes/analytics.py	py	4.2

`.env.example`

```
# SIEIS Environment Configuration
# Copy this file to .env and update values as needed

# =====
# Kafka Configuration
# =====
KAFKA_BROKER=redpanda:9092
KAFKA_TOPIC=sensor_readings

# =====
# InfluxDB 3.x Configuration (Real-time Hot Data, 30-day retention)
# =====
INFLUX_URL=http://influxdb3:8181
INFLUX_TOKEN=my-super-secret-token
INFLUX_DATABASE=sensor_data

# Note: InfluxDB 3.x uses DATABASE instead of ORG/BUCKET (v2.x concepts
removed)
# Data is automatically retained for 30 days (configured in docker-
compose.yml)

# =====
# MinIO Configuration (Historical Cold Data, Parquet Archive)
# =====
MINIO_ENDPOINT=minio:9000
MINIO_ACCESS_KEY=minioadmin
MINIO_SECRET_KEY=minioadmin123
MINIO_BUCKET=sieis-archive
MINIO_SECURE=false

# Note: For production, change credentials and set MINIO_SECURE=true

# =====
# Simulator Configuration
# =====
SPEED_FACTOR=100
DATA_PATH=/app/data/raw/data.txt
MOTE_LOCS_PATH=/app/data/raw/mote_locs.txt

# SPEED_FACTOR controls time compression:
#   100 = 100x faster (1 hour of data replays in 36 seconds)
#   1 = real-time (1:1 ratio)
#   1000 = 1000x faster (very fast for testing)

# =====
# Container vs Host Mode
# =====
# When running in Docker containers, services use internal DNS names (above)
```

```
# When running on host (python -m src.app...), use localhost:
#   KAFKA_BROKER=localhost:9092
#   INFLUX_URL=http://localhost:8181
#   MINIO_ENDPOINT=localhost:9000
```

Dockerfile.api

```
FROM python:3.11-slim
```

```
WORKDIR /app
```

```
COPY requirements.txt .
```

```
RUN pip install --no-cache-dir -r requirements.txt
```

```
COPY src/ ./src/
```

```
COPY .env.example .env
```

```
# Create models directory
```

```
RUN mkdir -p src/app/ml/models
```

```
CMD ["python", "-m", "src.app.api_server"]
```

Dockerfile.consumer

```
FROM python:3.11-slim
```

```
WORKDIR /app
```

```
# Install curl for healthchecks
```

```
RUN apt-get update && apt-get install -y curl && rm -rf /var/lib/apt/lists/*
```

```
# Copy requirements first for better layer caching
```

```
COPY requirements.txt .
```

```
# Install dependencies
```

```
RUN pip install --no-cache-dir -r requirements.txt
```

```
# Copy source code
```

```
COPY src/ ./src/
```

```
# Set Python path
```

```
ENV PYTHONUNBUFFERED=1
```

```
# Run consumer
```

```
CMD ["python", "-m", "src.app.consumer.main"]
```

Dockerfile.dashboard

```
FROM python:3.11-slim

WORKDIR /app

COPY requirements.txt .
RUN pip install --no-cache-dir -r requirements.txt

COPY src/ ./src/
COPY data/ ./data/
COPY .env.example .env

CMD ["streamlit", "run", "src/app/dashboard/app.py", "--server.port=8501", "--server.address=0.0.0.0", "--browser.gatherUsageStats=false"]
```

Dockerfile.scheduler

```
# — SIEIS Scheduler Container


---


# Runs APScheduler with two daily cron jobs:
#   Job A 00:00 UTC — remap incremental timestamps + restart simulator
#   Job B 02:00 UTC — retrain anomaly model + hot-reload FastAPI
#
# Requires /var/run/docker.sock mounted (to restart sieis-simulator).
#


---



FROM python:3.11-slim

WORKDIR /app

# Install system deps (curl for healthchecks)
RUN apt-get update && apt-get install -y --no-install-recommends \
    curl \
    && rm -rf /var/lib/apt/lists/*

# Copy and install Python dependencies
COPY requirements.txt .
RUN pip install --no-cache-dir -r requirements.txt

# Copy full source so scheduler can import src.app.* and scripts.*
COPY src/ ./src/
COPY scripts/ ./scripts/

# PYTHONPATH includes /app so both 'src.*' and 'scripts.*' resolve correctly
ENV PYTHONPATH=/app

# — Runtime defaults (override via docker-compose environment)


---


```

```
ENV SCHEDULER_TIMEZONE=UTC
ENV JOB_A_HOUR=0
ENV JOB_A_MINUTE=0
ENV JOB_B_HOUR=2
ENV JOB_B_MINUTE=0
ENV RUN_JOBS_ON_START=false
ENV SIMULATOR_CONTAINER=sieis-simulator
ENV API_RELOAD_URL=http://sieis-api:8000/api/v1/ml/model/reload
ENV INCR_DATA_PATH=/app/data/processed/incremental_data.txt

CMD ["python", "-m", "src.app.scheduler.main"]
```

Dockerfile.simulator

```
FROM python:3.11-slim

WORKDIR /app

# Copy requirements first for better layer caching
COPY requirements.txt .

# Install dependencies
RUN pip install --no-cache-dir -r requirements.txt

# Copy source code
COPY src/ ./src/

# Note: data files are mounted as volumes at runtime

# Set Python path
ENV PYTHONUNBUFFERED=1

# Run simulator
CMD ["python", "-m", "src.app.simulator.main"]
```

Documentation/QUICKSTART.md

SIEIS Quick Start

Run the full SIEIS stack (Kafka + Consumer + InfluxDB + MinIO + API + Dashboard + Scheduler) in about 10 minutes.

1. Prerequisites

- Docker Desktop running
- Python 3.11+ (optional, for utility scripts)
- Git

2. Required Data Files

Confirm these files exist:

- `data/raw/mote_locs.txt`
- `data/processed/incremental_data.txt`
- `data/processed/historical_data.txt`

If missing, regenerate with project scripts under `data/realtime_mapping/` and `scripts/`.

3. Start the Full Stack

From project root:

```
```powershell
docker-compose up --build -d
```
```

Check containers:

```
```powershell
docker ps --format "table {{.Names}}`t`{{.Status}}`t`{{.Ports}}"
```
```

Expected key containers:

- `sieis-redpanda`
- `sieis-console`
- `sieis-influxdb3`
- `sieis-minio`
- `sieis-minio-init`
- `sieis-simulator`
- `sieis-consumer`
- `sieis-api`
- `sieis-dashboard`
- `sieis-scheduler`

4. Verify Services Manually

4.1 Core health checks

```
```powershell
docker exec sieis-redpanda rpk cluster info
curl http://localhost:8086/health
curl http://localhost:9000/minio/health/live
curl http://localhost:8000/api/v1/health
```
```

4.2 Web consoles

- Redpanda Console: `http://localhost:8080`
- InfluxDB UI: `http://localhost:8086`
- MinIO Console: `http://localhost:9001`
- Dashboard: `http://localhost:8501`
- API docs: `http://localhost:8000/docs`

4.3 Login credentials

- InfluxDB
 - Username: `admin`
 - Password: `password123`
 - Org: `sieis`
 - Bucket: `sensor_data`
 - Token: `my-super-secret-token`
- MinIO
 - Username: `minioadmin`
 - Password: `minioadmin123`
 - Bucket: `sieis-archive`

5. Verify Streaming and Storage

Consume sample Kafka messages:

```
```powershell
docker exec sieis-redpanda rpk topic consume sensor_readings --num 5
```
```

Optional validation scripts:

```
```powershell
python scripts/verify_influxDb.py
python scripts/verify_minio_storage.py
```
```

6. ML Quick Check (Optional)

Train anomaly model manually:

```
```powershell
python scripts/train_model.py --source minio
```
```

Confirm model is loaded:

```
```powershell
curl http://localhost:8000/api/v1/ml/model/info
```
```

7. Scheduler Quick Check (Optional)

Run scheduler jobs immediately for smoke testing:

```
```powershell
docker compose down scheduler
RUN_JOBS_ON_START=true docker compose up scheduler
```
```

Default schedule (UTC):

- Job A `00:00`: restart simulator
- Job B `02:00`: retrain model and reload API model

8. Stop or Reset

Stop all services:

```
```powershell
docker-compose down
```
```

Stop and delete volumes (destructive reset):

```
```powershell
docker-compose down -v
```
```

9. Troubleshooting

- No dashboard data:
 - `docker logs sieis-consumer`
 - `docker logs sieis-simulator`
- API health degraded:
 - confirm InfluxDB reachable on `8086`
- MinIO empty:
 - check `sieis-minio-init` logs and bucket `sieis-archive`
- Scheduler not acting:
 - `docker logs sieis-scheduler --tail 200`

[README.md](#)

SIEIS: Smart Indoor Environmental Intelligence System

Abstract

The Smart Indoor Environmental Intelligence System (SIEIS) is a distributed data engineering and analytics platform for streaming environmental telemetry. The project integrates event-driven ingestion, dual-path persistence, API-based serving, dashboard-based observability, and scheduled model lifecycle operations. The objective is to support both low-latency operational monitoring and long-term analytical/ML workflows using a unified architecture.

1. System Architecture

1.1 Architectural Pattern

SIEIS follows a dual-write architecture:

- Hot path for near real-time monitoring and API queries (InfluxDB).
- Cold path for durable, partitioned, analytical storage (MinIO Parquet).

1.2 End-to-End Data Flow

```text

Processed Sensor File (incremental\_data.txt)

- > Simulator (threaded per mote)
- > Redpanda/Kafka topic: sensor\_readings
- > Consumer (batch polling + transformation)
- > InfluxDB (hot time-series store)
- > MinIO (cold Parquet archive)

Serving Layer:

- FastAPI queries InfluxDB and serves ML inference endpoints
- Streamlit dashboard consumes InfluxDB/MinIO/API outputs

Automation Layer:

- APScheduler service restarts simulator daily and retrains/reloads anomaly model

```

1.3 Runtime Topology (Docker Compose)

```mermaid

flowchart LR

```
A[Simulator] --> B[Redpanda/Kafka\nsensor_readings]
B --> C[Consumer]
C --> D[InfluxDB\nHot Storage]
C --> E[MinIO\nCold Storage]
D --> F[FastAPI]
E --> F
D --> G[Streamlit Dashboard]
E --> G
F --> G
H[Scheduler] --> A
H --> F
```

```

1.4 Architecture Diagram

![SIEIS Architecture] (Documentation/sieis-system-architecture.png)

2. Services and Port Addresses (Manual Verification)

The following table documents all externally exposed ports from the current `docker-compose.yml`.

Service	Container	Host Port(s)	Internal Address	Manual Verification
Redpanda (Kafka broker)	`sieis-redpanda`	`19092`, `9092`	`redpanda:9092`	`docker exec sieis-redpanda rpk cluster info`
Redpanda Console	`sieis-console`	`8080`	`redpanda-console:8080`	Open `http://localhost:8080`
InfluxDB	`sieis-influxdb3`	`8086`	`influxdb3:8086`	Open `http://localhost:8086/health`
MinIO API	`sieis-minio`	`9000`	`minio:9000`	Open `http://localhost:9000/minio/health/live`
MinIO Console	`sieis-minio`	`9001`	`minio:9001`	Open `http://localhost:9001`
FastAPI	`sieis-api`	`8000`	`sieis-api:8000`	Open `http://localhost:8000/api/v1/health`
Streamlit Dashboard	`sieis-dashboard`	`8501`	`sieis-dashboard:8501`	Open `http://localhost:8501`
Scheduler	`sieis-scheduler`	none (no host bind)	internal only	`docker logs sieis-scheduler`

Important note:

- Host-side Kafka clients should use `localhost:19092`.
- Inter-container Kafka traffic uses `redpanda:9092`.

2.1 Console Credentials and Login Details

Default local development credentials from the current compose configuration:

Service	Login URL	Username	Password	Additional Details
Redpanda Console	`http://localhost:8080`	Not required	Not required	Kafka UI only; no auth in local setup
InfluxDB UI	`http://localhost:8086`	`admin`	`password123`	Org: `sieis`, Bucket: `sensor_data`, Token: `my-super-secret-token`
MinIO Console	`http://localhost:9001`	`minioadmin`	`minioadmin123`	API endpoint: `http://localhost:9000`, Bucket: `sieis-archive`
FastAPI Docs	`http://localhost:8000/docs`	Not required	Not required	Health endpoint: `http://localhost:8000/api/v1/health`

Security note:

- These credentials are intended for local academic/development use only.
- Rotate and externalize secrets before any public or production deployment.

3. Local Prerequisites

3.1 Platform Requirements

- Operating System: Windows 10/11, Linux, or macOS
- Docker Desktop (or Docker Engine + Compose plugin)
- Python 3.11+

- Git

3.2 Python Dependencies

Install all Python packages from:

- `requirements.txt`

3.3 Required Data Files

The following files should exist before pipeline execution:

- `data/raw/mote_locs.txt`
- `data/processed/incremental_data.txt` (simulator input)
- `data/processed/historical_data.txt` (historical loading and ML utility workflows)

If processed files are absent or stale, regenerate using scripts in:

- `data/realtime_mapping/`
- `scripts/`

4. Procedure: Run the Project Locally

Step 1: Clone and enter the repository

```
```powershell
git clone <repository-url>
cd SIEIS
```
```

Step 2: (Optional but recommended) Create Python virtual environment

```
```powershell
python -m venv venv
.\venv\Scripts\Activate.ps1
pip install -r requirements.txt
```
```

Step 3: Start full system with Docker Compose

```
```powershell
docker-compose up --build -d
```
```

Step 4: Validate container startup

```
```powershell
docker ps --format "table {{.Names}}\t{{.Status}}\t{{.Ports}}"
```
```

Expected containers include:

- `sieis-redpanda`

- `sieis-console`
- `sieis-influxdb3`
- `sieis-minio`
- `sieis-minio-init`
- `sieis-simulator`
- `sieis-consumer`
- `sieis-api`
- `sieis-dashboard`
- `sieis-scheduler`

Step 5: Manual service verification

1. Verify Kafka broker

```
```powershell
docker exec sieis-redpanda rpk cluster info
```
```

2. Verify InfluxDB

```
```powershell
curl http://localhost:8086/health
```
```

3. Verify MinIO

```
```powershell
curl http://localhost:9000/minio/health/live
```
```

4. Verify API

```
```powershell
curl http://localhost:8000/api/v1/health
```
```

5. Verify message flow

```
```powershell
docker exec sieis-redpanda rpk topic consume sensor_readings --num 5
```
```

6. Open web interfaces

- Redpanda Console: `http://localhost:8080`
- InfluxDB UI: `http://localhost:8086`
- MinIO Console: `http://localhost:9001`
- Dashboard: `http://localhost:8501`

Step 6: Shut down environment

```
```powershell
docker-compose down
```
```

Remove volumes (destructive) when complete reset is required:

```
```powershell
docker-compose down -v
```
```

5. Local Development Execution (Without Full Compose)

For targeted debugging, services can be run directly from Python entry points:

```
```powershell
python -m src.app.simulator.main
python -m src.app.consumer.main
python -m src.app.api_server
python scripts/run_dashboard.py
```
```

Use host-compatible environment values in `.env` (for example, `localhost` endpoints) when running services outside containers.

6. Project Folder Structure

```
```text
SIEIS/
|-- docker-compose.yml
|-- requirements.txt
|-- README.md
|-- QUICKSTART.md
|-- .env.example
|-- data/
| |-- raw/
| | |-- mote_locs.txt
| |-- processed/
| | |-- historical_data.txt
| | |-- incremental_data.txt
| | |-- realtime_data.txt
| |-- realtime_mapping/
| |-- transform_to_realtime.py
| |-- preview_mapping.py
| |-- validate_output.py
|-- src/
| |-- app/
| |-- config.py
| |-- api_server.py
| |-- api/
```

```

| | |-- main.py
| | |-- schemas.py
| | `-- routes/
| | |-- sensors.py
| | |-- analytics.py
| | `-- ml.py
| |-- simulator/
| | |-- main.py
| | |-- orchestrator.py
| | |-- emitter.py
| | |-- producer.py
| | `-- data_loader.py
| |-- consumer/
| | |-- main.py
| | |-- kafka_consumer.py
| | |-- influx_writer.py
| | `-- parquet_writer.py
| |-- dashboard/
| | |-- app.py
| | `-- pages/
| | |-- 1_Realtime_Monitor.py
| | |-- 2_Historical_Analysis.py
| | `-- 3_Anomaly_Detection.py
| |-- ml/
| | |-- detector.py
| | |-- models/
| | | `-- model_registry.json
| | `-- preprocessing/
| | `-- data_prep.py
| `-- scheduler/
| |-- main.py
| `-- jobs.py
|-- scripts/
| |-- load_historical_data.py
| |-- split_dataset.py
| |-- remap_timestamps.py
| |-- train_model.py
| |-- retrain_model.py
| `-- verify_*.py
|-- tests/
| |-- test_e2e_pipeline.py
| |-- test_full_pipeline.py
| |-- test_data_loader.py
| `-- test_container_*.py
`-- Documentation/
 |-- ARCHITECTURE.md
 |-- DEPLOYMENT.md
 |-- ML_MODELS.md
 `-- PROJECT_CONTEXT.md
...

```



Structure note:

- `src/app/` contains runtime services and application logic.
- `scripts/` contains operational utilities (data prep, loading, verification, model workflows).
- `data/processed/` stores the categorized datasets used to emulate archival plus incremental IoT flow.

## ## 7. Troubleshooting Quick Reference

- If API is healthy but dashboard is empty, verify Kafka consumption and Influx writes via `docker logs sieis-consumer`.
- If simulator produces no current records, regenerate mapped data in `data/realtime\_mapping` and restart `sieis-simulator`.
- If MinIO appears empty, verify bucket initialization from `sieis-minio-init` logs.
- If model endpoints use fallback mode, train/retrain a model and call `POST /api/v1/ml/model/reload`.

## ## 8. Data Source and Temporal Categorization for Real-Time IoT Emulation

### ### 8.1 Source Dataset

SIEIS uses the Intel Lab sensor dataset (historical indoor environmental telemetry), where each record includes:

- date and time
- epoch/sample index
- mote identifier
- temperature, humidity, light, and voltage

The original dataset is historical (early 2000s), so it does not naturally behave as a present-day live stream.

### ### 8.2 Preprocessing and Time Remapping

To make legacy telemetry operationally useful for current-time demonstrations, the pipeline adds a derived field:

- `updated\_timestamp`: a remapped timestamp aligned to the present calendar window

This preserves the original intra-series behavior (ordering and temporal spacing) while shifting records into an interpretable modern timeline.

### ### 8.3 Historical vs Incremental Categorization

The project organizes processed data into three artifacts under `data/processed/`:

1. `realtime\_data.txt`
  - Full remapped dataset (complete timeline representation).

## 2. `historical\_data.txt`

- Archive-oriented partition (typically the older 80% block after chronological split).
- Used for backfill, long-horizon analysis, and model training workflows.
- Loaded through historical loaders into MinIO Parquet (and optionally InfluxDB as needed).

## 3. `incremental\_data.txt`

- Stream-oriented partition (typically the newer 20% block).
- Used by the simulator as the active feed for Kafka publication.
- Represents the "arriving" portion of data to mimic ongoing sensor production.

### ### 8.4 Why This Mimics a Real IoT Scenario

This design reproduces the operational structure of production IoT systems:

- historical corpus retained in a cold analytical tier (MinIO/Parquet),
- incremental events continuously emitted through a message bus (Kafka/Redpanda),
- low-latency monitoring served from hot time-series storage (InfluxDB),
- periodic model retraining on accumulated historical data.

Consequently, SIEIS can emulate:

- streaming ingestion dynamics,
- online dashboard/API consumption,
- offline analytical and ML lifecycle tasks,
- and the natural separation between "already observed" and "newly arriving" telemetry.

### ## 9. ML Training, Retraining, and Scheduler Automation

#### ### 9.1 ML Component Overview

SIEIS uses an Isolation Forest anomaly detector implemented in:

- `src/app/ml/detector.py`
- `src/app/ml/preprocessing/data\_prep.py`

Model artifacts are stored in:

- `src/app/ml/models/anomaly\_detector\_\*.pkl`

The active model pointer is tracked in:

- `src/app/ml/models/model\_registry.json`

FastAPI serves inference and model-management endpoints:

- `POST /api/v1/ml/predict/anomaly`
- `GET /api/v1/ml/model/info`
- `POST /api/v1/ml/model/reload`

#### ### 9.2 Initial Model Training (Manual)

Run this after historical data is available (local file or MinIO Parquet):

```
```powershell
python scripts/train_model.py --source minio
```
```

Useful variants:

```
```powershell
python scripts/train_model.py --source local --max-rows 200000
python scripts/train_model.py --source minio --contamination 0.03
python scripts/train_model.py --source minio --tag baseline_v1
```
```

What the script does:

1. Loads training data (`minio` or `local` source)
2. Prepares features (`temperature`, `humidity`, `light`, `voltage`, `hour`, `day\_of\_week`)
3. Trains Isolation Forest
4. Saves model artifact and updates `model\_registry.json`
5. Attempts API model reload and sends a test prediction

### ### 9.3 Incremental Retraining (Manual)

Use retraining when new data arrives and you want to refresh model behavior:

```
```powershell
python scripts/retrain_model.py --source minio --days 30
```
```

Alternative source:

```
```powershell
python scripts/retrain_model.py --source influxdb --days 7
```
```

After retraining, reload the API model if not automatically reloaded:

```
```powershell
curl -X POST http://localhost:8000/api/v1/ml/model/reload
```
```

### ### 9.4 Scheduler-Based Daily Automation

The `scheduler` container (`src/app/scheduler/main.py`) runs two UTC cron jobs:

- Job A (`00:00 UTC`): `remap\_and\_restart`
  - Current implementation restarts `sieis-simulator`.

- This ensures the simulator process starts a fresh emission cycle.
- Job B (`02:00 UTC`): `retrain\_and\_reload`
  - Loads recent data from MinIO Parquet.
  - Trains a new Isolation Forest model.
  - Saves a timestamped model with `scheduled` tag.
  - Calls API reload endpoint so inference uses the latest model without container restart.

Immediate smoke test (run both jobs on startup):

```
```powershell
docker compose up -d scheduler
docker compose down scheduler
RUN_JOBS_ON_START=true docker compose up scheduler
```
```

Scheduler runtime configuration is controlled by environment variables in `docker-compose.yml`, including:

- `SCHEDULER\_TIMEZONE`
- `JOB\_A\_HOUR`, `JOB\_A\_MINUTE`
- `JOB\_B\_HOUR`, `JOB\_B\_MINUTE`
- `RUN\_JOBS\_ON\_START`

### ### 9.5 Verification After Training/Retraining

1. Confirm API sees the model:

```
```powershell
curl http://localhost:8000/api/v1/ml/model/info
```
```

2. Confirm health includes model-loaded state:

```
```powershell
curl http://localhost:8000/api/v1/health
```
```

3. Confirm scheduler execution logs:

```
```powershell
docker logs sieis-scheduler --tail 200
```
```

4. Confirm new model files are created:

```
```powershell
Get-ChildItem src\app\ml\models\anomaly_detector_*.pkl | Sort-Object
LastWriteTime -Descending | Select-Object -First 5
```
```

...

data/.checkpoint\_historical.json

```
{
 "last_line": 0,
 "timestamp": "2026-02-25T11:44:52.529371",
 "stats": {
 "total_read": 0,
 "influx_written": 0,
 "minio_written": 0,
 "skipped": 0,
 "errors": 0
 },
 "note": "Reset by clear_storage.py"
}
```

data/raw/mote\_locs.txt

```
1 21.5 23
2 24.5 20
3 19.5 19
4 22.5 15
5 24.5 12
6 19.5 12
7 22.5 8
8 24.5 4
9 21.5 2
10 19.5 5
11 16.5 3
12 13.5 1
13 12.5 5
14 8.5 6
15 5.5 3
16 1.5 2
17 1.5 8
18 5.5 10
19 3.5 13
20 0.5 17
21 4.5 18
22 1.5 23
23 6 24
24 1.5 30
25 4.5 30
26 7.5 31
27 8.5 26
28 10.5 31
29 12.5 26
30 13.5 31
```

```
31 15.5 28
32 17.5 31
33 19.5 26
34 21.5 30
35 24.5 27
36 26.5 31
37 27.5 26
38 30.5 31
39 30.5 26
40 33.5 28
41 36.5 30
42 39.5 30
43 35.5 24
44 40.5 22
45 37.5 19
46 34.5 16
47 39.5 14
48 35.5 10
49 39.5 6
50 38.5 1
51 35.5 4
52 31.5 6
53 28.5 5
54 26.5 2
```

## [data/realtime\\_mapping/README.md](#)

# Real-Time Data Mapping Scripts

This folder contains scripts to transform historical sensor data (2004) into real-time relevant timestamps (2026).

## Strategy

Maps the date range proportionally with an 80% anchor point:

- The date at 80% of the original range → Today's date
- All other dates map proportionally backwards
- Future dates are filtered out

**Example:**

- Original: 2004-02-28 to 2004-05-01 (63 days)
- 80% point: 2004-04-18 → 2026-02-16 (today)
- 2004-02-28 → 2026-01-28 (19 days ago)
- 2004-04-17 → 2026-02-15 (yesterday)

## Scripts

### 1. `preview\_mapping.py`

Preview the date mapping before transformation.

```
```bash
python data/realtime_mapping/preview_mapping.py
```
```

**\*\*Output:\*\*** Shows how each original date will map to new dates.

### 2. `transform\_to\_realtime.py`

Main transformation script that creates the real-time data file.

```
```bash
python data/realtime_mapping/transform_to_realtime.py
```
```

**\*\*Input:\*\*** `data/data.csv` (original)

**\*\*Output:\*\***

- `data/data\_realtime.csv` (CSV format with headers)

- `data/processed/realtime\_data.txt` (Space-delimited format, no headers)

### 3. `validate\_output.py`

Validates the transformed data file.

```
```bash
python data/realtime_mapping/validate_output.py
```
```

**\*\*Checks:\*\***

- Date distribution
- Timestamp consistency
- No future dates
- Record counts

## Usage Workflow

```
```bash
# Step 1: Preview the mapping (optional)
python data/realtime_mapping/preview_mapping.py

# Step 2: Transform the data
python data/realtime_mapping/transform_to_realtime.py

# Step 3: Validate output (optional)
python data/realtime_mapping/validate_output.py

# Step 4: Update docker-compose.yml
# Change simulator volume from:
#   - ./data/data.csv:/data/data.csv:ro
# To:
#   - ./data/data_realtime.csv:/data/data.csv:ro
```
```

```

Step 5: Restart containers
docker-compose restart simulator consumer

Step 6: Run tests
python tests/test_full_pipeline.py
```

## Output Formats

### CSV Format (`data/data_realtime.csv`)
Standard CSV with headers:
```csv
date,time,epoch,moteid,temperature,humidity,light,voltage,updated_timestamp
2004-02-28,00:59:16.02785,3,1,19.9884,37.0933,45.08,2.69964,2026-01-
17T00:59:16.027850
```

**Columns:**
- `date`: Original date from 2004
- `time`: Original time of day
- `epoch`: Epoch/batch number
- `moteid`: Sensor/mote identifier
- `temperature`: Temperature reading (°C)
- `humidity`: Humidity reading (%)
- `light`: Light reading (lux)
- `voltage`: Battery voltage (V)
- `updated_timestamp`: **New** - Mapped timestamp in 2026 (ISO 8601 format)

### TXT Format (`data/processed/realtime_data.txt`)
Space-delimited, no headers (for legacy systems):
```
2004-02-28 00:59:16.02785 3 1 19.9884 37.0933 45.08 2.69964 2026-01-
17T00:59:16.027850
```

**Column order:**
1. date (original)
2. time (original)
3. epoch
4. moteid
5. temperature
6. humidity
7. light
8. voltage
9. updated_timestamp (mapped to 2026)

## Transformation Details

### Date Mapping Algorithm

```


1. ****Analyze original dataset:****
 - Find min date, max date, total span
 - Calculate 80% point (day 50 of 63 = 2004-04-18)
2. ****Create proportional mapping:****
 - Anchor: Date at 80% → Today
 - All dates map relative to this anchor
 - Preserves time-of-day from original timestamps
3. ****Filter future dates:****
 - Any timestamp > now() is excluded
 - Prevents InfluxDB rejection (can't write future data)

Example Transformation

...

Original Range: 2004-02-28 to 2004-05-01 (63 days)

Current Date: 2026-02-16

Mapping:

| | | | | | |
|----|------------|---|------------|-------------------------|----------|
| └─ | 2004-02-28 | → | 2026-01-28 | (19 days ago) | ← Oldest |
| └─ | 2004-03-15 | → | 2026-02-02 | (14 days ago) | |
| └─ | 2004-04-17 | → | 2026-02-15 | (1 day ago) | |
| └─ | 2004-04-18 | → | 2026-02-16 | (TODAY) ★ 80% anchor | |
| └─ | 2004-04-19 | → | 2026-02-17 | (FUTURE - filtered out) | |
| └─ | 2004-05-01 | → | 2026-03-01 | (FUTURE - filtered out) | ← Newest |

...

Statistics from Last Run

...

Original Dataset:

Records: 2,313,682

Date Range: 2004-02-28 to 2004-05-01

Unique Dates: 63

80% Point: 2004-04-18

Transformation Results:

Mapped Records: 1,823,589 (78.8%)

Filtered (future): 490,093 (21.2%)

Distribution:

Today: 856,234 records (47.0%)

Historical (1-30 days ago): 967,355 records (53.0%)

...

File Locations

...

```

data/
├─ data.csv                # Original (untouched)
├─ data_realtime.csv       # Transformed CSV
└─ processed/
    └─ realtime_data.txt   # Transformed TXT (space-delimited)

data/realtime_mapping/
├─ README.md               # This file
├─ preview_mapping.py      # Preview script
├─ transform_to_realtime.py # Main transformation
└─ validate_output.py      # Validation script
...

## Integration with SIEIS Pipeline

### Before Transformation
...
CSV (2004 dates) → Simulator → Kafka → Consumer → InfluxDB ❌ (rejected)
                                                         → MinIO ✅ (works)
...

**Issue:** InfluxDB rejects data outside retention window (30 days default)

### After Transformation
...
CSV (2026 dates) → Simulator → Kafka → Consumer → InfluxDB ✅ (accepted)
                                                         → MinIO ✅ (works)
...

**Result:** Both hot path (InfluxDB) and cold path (MinIO) work correctly

## Notes

- ✅ Preserves time-of-day from original timestamps for realism
- ✅ Filters out any records that would map to future dates
- ✅ Safe to re-run multiple times (overwrites output files)
- ✅ Does not modify original `data/data.csv`
- ✅ Creates `data/processed/` directory if it doesn't exist
- ✅ Both CSV and TXT formats generated automatically
- ⚠️ TXT format has NO headers (by design for legacy compatibility)
- ⚠️ Re-run transformation when testing on different days (dates shift)

## Troubleshooting

### Problem: "No data in InfluxDB"
**Solution:** Ensure you're using `data_realtime.csv` in docker-compose.yml

### Problem: "All data filtered out"

```

****Cause:**** Running script on a different date than dataset was created
****Solution:**** Re-run ``transform_to_realtime.py`` to regenerate with current date

Problem: "Future timestamp rejected"

****Cause:**** System clock incorrect or transformation bug

****Solution:****

1. Check system date: ``date`` (Linux/Mac) or ``Get-Date`` (PowerShell)
2. Validate output: ``python data/realtime_mapping/validate_output.py``

Problem: "CSV has headers but TXT doesn't"

****Expected:**** This is by design. Use CSV for imports, TXT for legacy systems.

Version History

- ****v1.0.0**** (2026-02-16): Initial release with 80% proportional mapping
- ****v1.1.0**** (2026-02-16): Added TXT format output to ``data/processed/``

Related Documentation

- [SIEIS Architecture](../..//Documentation/ARCHITECTURE.md)
- [Retention Policy Guide](../..//Documentation/RETENTION_POLICY_GUIDE.md)
- [Testing Guide](../..//tests/README.md)

`data/realtime_mapping/preview_mapping.py`

"""

Preview how dates will be mapped without creating the full file.
Quick way to verify the mapping logic before processing.

"""

```
import csv
```

```
from datetime import datetime, timedelta
```

```
from pathlib import Path
```

```
import logging
```

```
logging.basicConfig(level=logging.INFO, format='%(message)s')
```

```
logger = logging.getLogger(__name__)
```

```
def main():
```

```
    # Navigate to project root from data/realtime_mapping/
```

```
    script_dir = Path(__file__).parent
```

```
    project_root = script_dir.parent.parent
```

```
    input_file = project_root / "data" / "data.csv"
```

```
    if not input_file.exists():
```

```
        logger.error(f"❌ File not found: {input_file}")
```

```
        logger.info(f"\nSearching for data files...")
```

```

# Try alternative locations
alt_locations = [
    project_root / "data" / "raw" / "data.txt",
    project_root / "data" / "data.txt",
]

for alt_file in alt_locations:
    if alt_file.exists():
        logger.info(f"Found: {alt_file}")
        input_file = alt_file
        break
    else:
        logger.error(f"\n❌ No data file found. Please ensure data.csv
exists in:")
        logger.error(f"    {project_root / 'data' / 'data.csv'}")
        return 1

logger.info("🔍 Previewing Date Mapping")
logger.info(f"📁 Input: {input_file.relative_to(project_root)}\n")

# Read all unique dates
dates = []
with open(input_file, 'r') as f:
    # Check if space-delimited or CSV
    first_line = f.readline()
    f.seek(0)

    if ',' in first_line:
        # CSV format
        reader = csv.DictReader(f)
        for row in reader:
            try:
                ts = datetime.strptime(row['date'], "%Y-%m-%d
%H:%M:%S.%f")
            except (ValueError, KeyError):
                continue
            dates.append(ts.date())
    else:
        # Space-delimited format: date time mote_id epoch temp humidity
        light voltage
        for line in f:
            parts = line.strip().split()
            if len(parts) >= 2:
                try:
                    timestamp_str = f"{parts[0]} {parts[1]}"
                    ts = datetime.strptime(timestamp_str, "%Y-%m-%d
%H:%M:%S.%f")
                    dates.append(ts.date())

```

```

        except ValueError:
            continue

if not dates:
    logger.error("❌ No valid dates found in file")
    return 1

unique_dates = sorted(set(dates))
min_date = unique_dates[0]
max_date = unique_dates[-1]
total_days = (max_date - min_date).days + 1

# Calculate 80% point
days_to_80 = int(total_days * 0.80)
date_at_80 = min_date + timedelta(days=days_to_80)

today = datetime.now().date()

logger.info(f"📊 Original Dataset:")
logger.info(f"    Date Range: {min_date} to {max_date}")
logger.info(f"    Total Days: {total_days}")
logger.info(f"    Unique Dates: {len(unique_dates)}")
logger.info(f"    80% Point: {date_at_80} (day {days_to_80})")

logger.info(f"\n🗺️ Mapping Strategy:")
logger.info(f"    {date_at_80} → {today} (TODAY)")
logger.info(f"    All dates map relative to this anchor")

logger.info(f"\n📅 Complete Date Mapping:\n")
logger.info(f"    {'Original Date':<15} → {'New Date':<15}"
{'Description'})
logger.info(f"    {'-'*15}    {'-'*15} {'-'*30}")

# Show all mappings
for orig_date in unique_dates:
    days_offset = (orig_date - date_at_80).days
    new_date = today + timedelta(days=days_offset)

    # Description
    diff = (new_date - today).days
    if diff == 0:
        desc = "TODAY 🌟"
    elif diff < 0:
        desc = f"{abs(diff)} days ago"
    else:
        desc = f"{diff} days in FUTURE (filtered)"

    # Highlight special dates
    marker = ""

```

```

        if orig_date == min_date:
            marker = " ← oldest"
        elif orig_date == max_date:
            marker = " ← newest"
        elif orig_date == date_at_80:
            marker = " ← 80% anchor"

    logger.info(f"    {orig_date}    →    {new_date}    {desc}{marker}")

# Count records per date
logger.info(f"\n📊 Records per Original Date:")
date_counts = {}
for d in dates:
    date_counts[d] = date_counts.get(d, 0) + 1

for date in sorted(date_counts.keys())[:10]:
    count = date_counts[date]
    logger.info(f"    {date}: {count:,} records")

if len(date_counts) > 10:
    logger.info(f"    ... ({len(date_counts) - 10} more dates)")

# Calculate expected distribution
future_count = sum(1 for d in dates if (today + timedelta(days=(d -
date_at_80).days)) > today)
past_count = sum(1 for d in dates if (today + timedelta(days=(d -
date_at_80).days)) <= today)

logger.info(f"\n✅ Expected Output:")
logger.info(f"    Total records: {len(dates):,}")
logger.info(f"    Will be kept: {past_count:,}
({past_count/len(dates)*100:.1f}%)")
logger.info(f"    Will be filtered (future): {future_count:,}
({future_count/len(dates)*100:.1f}%)")

logger.info(f"\n✅ Ready to transform? Run:")
logger.info(f"    python data/realtime_mapping/transform_to_realtime.py")

return 0

if __name__ == "__main__":
    exit(main())

```

[data/realtime_mapping/transform_to_realtime.py](#)

"""

Split raw sensor data into Historical and Incremental datasets with time mapping.

Logic:

1. Load raw data from data/raw/data.txt
 2. Sort by timestamp (ascending)
 3. Determine 80% split point based on time range
 4. Split data:
 - Historical: <= 80% date (Mapped so last date = TODAY, rest backwards)
 - Incremental: > 80% date (Mapped so first date = TOMORROW, rest forwards)
- """

```
import logging
import pandas as pd
from datetime import datetime, timedelta
from pathlib import Path
import sys

# Configure logging
logging.basicConfig(
    level=logging.INFO,
    format='%(asctime)s - %(levelname)s - %(message)s'
)
logger = logging.getLogger(__name__)

def parse_raw_data(file_path):
    """
    Read raw intel lab data.
    Format is space separated: date time epoch moteid temp humidity light
    voltage
    """
    logger.info(f"📖 Reading raw data from {file_path}...")

    # Define column names based on Intel Lab dataset format
    col_names = ['date', 'time', 'epoch', 'moteid', 'temperature', 'humidity',
                  'light', 'voltage']

    try:
        # Read using pandas for efficiency
        df = pd.read_csv(
            file_path,
            sep=r'\s+',
            header=None,
            names=col_names,
            on_bad_lines='skip',
            engine='python' # Python engine is more robust for variable
whitespace
        )

        # Create a single timestamp column for sorting
```

```

df['original_ts'] = pd.to_datetime(
    df['date'] + ' ' + df['time'],
    format='%Y-%m-%d %H:%M:%S.%f',
    errors='coerce'
)

# Drop invalid rows
initial_len = len(df)
df = df.dropna(subset=['original_ts'])
if len(df) < initial_len:
    logger.warning(f"⚠️ Dropped {initial_len - len(df)} rows with
invalid timestamps")

# Sort by timestamp
logger.info(f"⚡ Sorting data by timestamp...")
df = df.sort_values('original_ts')

return df

except Exception as e:
    logger.error(f"❌ Failed to parse data: {e}")
    raise

def calculate_split_date(df):
    """Find the date that represents the 80% mark of the time range."""
    min_date = df['original_ts'].min()
    max_date = df['original_ts'].max()

    total_duration = max_date - min_date
    split_offset = total_duration * 0.8
    split_date = min_date + split_offset

    logger.info(f"📅 Date Range: {min_date} to {max_date}")
    logger.info(f"🕒 Duration: {total_duration}")
    logger.info(f"✂️ 80% Split Timestamp: {split_date}")

    return split_date

def map_timestamps(df, target_anchor_date, is_historical=True):
    """
    Map timestamps to new range.

    Args:
        df: DataFrame slice
        target_anchor_date: The target date to map the specific anchor point
to.
        is_historical:

```



```

        If True: Map MAX date in df -> target_anchor_date (Today).
        If False: Map MIN date in df -> target_anchor_date (Tomorrow).
    """
    if df.empty:
        return df

    # Calculate offsets
    if is_historical:
        # For historical: Last record = Today (aligned at same time of day as
original max)
        anchor_original = df['original_ts'].max()
        time_offset = target_anchor_date - anchor_original
        logger.info(f"Mapping Historical: Max Original {anchor_original} ->
Target {target_anchor_date}")
    else:
        # For incremental: First record = Tomorrow (Start of day)
        anchor_original = df['original_ts'].min()
        time_offset = target_anchor_date - anchor_original
        logger.info(f"Mapping Incremental: Min Original {anchor_original} ->
Target {target_anchor_date}")

    # Apply offset
    df = df.copy()
    df['updated_timestamp'] = df['original_ts'] + time_offset

    return df

def save_processed_file(df, output_path):
    """Save the processed DataFrame to CSV/TXT format expected by
simulator."""
    if df.empty:
        logger.warning(f"⚠ Skipping empty dataframe save to {output_path}")
        return

    logger.info(f"💾 Saving {len(df):,} records to {output_path}...")

    # NEW LOGIC: Format updated_timestamp to be ISO format (no spaces) to
prevent parse errors downstream
    # pd.to_csv default string representation might include spaces.
    if 'updated_timestamp' in df.columns:
        # We need to operate on a copy or modify safely.
        # Since this function is the last step for these DFs, modifying in
place or copy is fine.
        df = df.copy()
        df['updated_timestamp'] = df['updated_timestamp'].dt.strftime('%Y-%m-
%dT%H:%M:%S.%f')

    # Ensure directory exists

```

```

output_path.parent.mkdir(parents=True, exist_ok=True)

# Columns to save: original columns + updated_timestamp
# Note: Simulator might expect specific column order. Usually it reads
columns by name if header/names provided in pd.read_csv
# Original raw data has no header.
# Our data_loader.py reads: 'date', 'time', 'epoch', 'moteid',
'temperature', 'humidity', 'light', 'voltage', 'updated_timestamp'

cols_to_save = ['date', 'time', 'epoch', 'moteid', 'temperature',
'humidity', 'light', 'voltage', 'updated_timestamp']

# Use space separator as per original format, no header, index=False
import csv as csv_module
df.to_csv(
    output_path,
    sep=' ',
    index=False,
    columns=cols_to_save,
    header=False,
    quoting=csv_module.QUOTE_NONE,
    escapechar=' '
)

def main():
    script_dir = Path(__file__).parent
    project_root = script_dir.parent.parent

    # Check possible input file locations
    possible_inputs = [
        project_root / "data" / "raw" / "data.txt",
        project_root / "data" / "data.txt"
    ]

    raw_file = None
    for p in possible_inputs:
        if p.exists():
            raw_file = p
            break

    if not raw_file:
        logger.error("✗ Raw data file not found in data/raw/data.txt")
        return 1

    # 1. Load and Sort
    try:
        df = parse_raw_data(raw_file)
    except Exception as e:

```

```

        logger.error(f"Failed to load data: {e}")
        return 1

# 2. Determine Split
split_timestamp = calculate_split_date(df)

# 3. Split Dataframes
mask_historical = df['original_ts'] <= split_timestamp
df_hist = df[mask_historical].copy()
df_incr = df[~mask_historical].copy()

logger.info(f"📊 Split Stats:")
logger.info(f"    Historical Records: {len(df_hist):,} (Max Orig Date: {df_hist['original_ts'].max():,})")
logger.info(f"    Incremental Records: {len(df_incr):,} (Min Orig Date: {df_incr['original_ts'].min() if not df_incr.empty else 'N/A':,})")

# 4. Map Dates
now = datetime.now()
today = now.replace(hour=0, minute=0, second=0, microsecond=0)

# Calculate target anchor for Historical (End of Today or Now?)
# "80th% percent date is mapped to today" implies the end of the
historical range is today.
# Let's map it to NOW.
target_hist_anchor = now

df_hist_mapped = map_timestamps(df_hist, target_hist_anchor,
is_historical=True)

# Calculate target anchor for Incremental
# User requested: "transmit only current day data from incremental data
file"
# To facilitate this, we start the incremental data TODAY (now + 1 minute)
# instead of tomorrow, so the simulator picks it up immediately.
target_incr_start = now + timedelta(minutes=1)

df_incr_mapped = map_timestamps(df_incr, target_incr_start,
is_historical=False)

# 5. Save Files
output_dir = project_root / "data" / "processed"

save_processed_file(df_hist_mapped, output_dir / "historical_data.txt")
save_processed_file(df_incr_mapped, output_dir / "incremental_data.txt")

# 6. Create 'realtime_data.txt' for simulator backward compatibility
# The simulator currently points to data/processed/realtime_data.txt
# We should probably combine them or just use historical?

```

```

    # Request implies splitting, maybe for a staged simulation (load
    historical, then stream incremental).
    # For now, I'll save the historical data as 'realtime_data.txt' too so the
    current simulator has something relevant to play.
    # OR better: save the concatenation of both mapped datasets to
    realtime_data.txt so the simulator can iterate through all.
    # Since mapped timestamps are contiguous (Today -> Tomorrow...), let's
    save the Combined set.

    logger.info("🔗 Creating combined realtime_data.txt for simulator
    compatibility...")
    df_combined = pd.concat([df_hist_mapped, df_incr_mapped])
    save_processed_file(df_combined, output_dir / "realtime_data.txt")

    logger.info("✅ Data transformation complete.")
    logger.info(f"    Output 1: {output_dir / 'historical_data.txt'}")
    logger.info(f"    Output 2: {output_dir / 'incremental_data.txt'}")
    logger.info(f"    Output 3: {output_dir / 'realtime_data.txt'} (Combined)")
    return 0

if __name__ == "__main__":
    sys.exit(main())

```

data/realtime_mapping/validate_output.py

```

"""
Validate the transformed real-time data file.
Checks for data integrity, distribution, and timestamp consistency.
"""

import csv
from datetime import datetime
from pathlib import Path
from collections import defaultdict
import logging

logging.basicConfig(level=logging.INFO, format='%(message)s')
logger = logging.getLogger(__name__)

def validate_output_file(output_file):
    """Perform comprehensive validation of the transformed data."""

    logger.info("🔍 Validating Transformed Data File")
    logger.info(f"📁 File: {output_file}\n")

    if not output_file.exists():
        logger.error(f"❌ File not found: {output_file}")
        logger.info("\nRun transformation first:")

```

```

        logger.info("    python
data/realtime_mapping/transform_to_realtime.py")
        return False

    # Validation metrics
    total_records = 0
    today_count = 0
    historical_count = 0
    future_count = 0

    mote_ids = set()
    date_distribution = defaultdict(int)
    fields_summary = defaultdict(lambda: {'count': 0, 'sum': 0, 'min':
float('inf'), 'max': float('-inf')})

    sample_records = []
    errors = []

    today = datetime.now().date()
    now = datetime.now()

    # Read and analyze
    logger.info("🇮🇹 Analyzing data...")

    with open(output_file, 'r') as f:
        reader = csv.DictReader(f)

        # Verify required columns
        required_cols = ['date', 'updated_timestamp', 'mote_id',
'temperature', 'humidity']
        missing_cols = [col for col in required_cols if col not in
reader.fieldnames]

        if missing_cols:
            logger.error(f"❌ Missing required columns: {missing_cols}")
            return False

        for i, row in enumerate(reader):
            total_records += 1

            try:
                # Parse timestamps
                updated_ts = datetime.fromisoformat(row['updated_timestamp'])
                updated_date = updated_ts.date()

                # Check for future dates (should not exist)
                if updated_ts > now:
                    future_count += 1
                    if len(errors) < 5:

```

```

        errors.append(f"Row {i+1}: Future timestamp
{updated_ts.isoformat()}")

    # Categorize by date
    if updated_date == today:
        today_count += 1
    else:
        historical_count += 1

    # Track distribution
    date_key = updated_date.isoformat()
    date_distribution[date_key] += 1

    # Track mote IDs
    mote_ids.add(row['mote_id'])

    # Analyze sensor fields
    for field in ['temperature', 'humidity', 'light', 'voltage']:
        try:
            value = float(row[field])
            fields_summary[field]['count'] += 1
            fields_summary[field]['sum'] += value
            fields_summary[field]['min'] =
min(fields_summary[field]['min'], value)
            fields_summary[field]['max'] =
max(fields_summary[field]['max'], value)
        except (ValueError, KeyError):
            pass

    # Collect samples
    if i < 3 or (i % 50000 == 0 and len(sample_records) < 10):
        sample_records.append({
            'mote_id': row['mote_id'],
            'original': row['date'],
            'updated': row['updated_timestamp'],
            'temp': row.get('temperature', 'N/A')
        })

    except Exception as e:
        if len(errors) < 5:
            errors.append(f"Row {i+1}: Parse error - {e}")

    # Progress
    if total_records % 50000 == 0:
        logger.info(f"    Processed {total_records:,} records...")

    # Print results
    logger.info(f"\n{'='*80}")
    logger.info("VALIDATION RESULTS")
    logger.info(f"{'='*80}\n")

```

```

# Basic stats
logger.info("📊 Record Counts:")
logger.info(f"    Total records: {total_records:,}")
logger.info(f"    Today's data: {today_count:,}
({today_count/total_records*100:.1f}%)")
logger.info(f"    Historical: {historical_count:,}
({historical_count/total_records*100:.1f}%)")
logger.info(f"    Unique notes: {len(note_ids)}")

# Check for errors
if future_count > 0:
    logger.error(f"\n❌ VALIDATION FAILED!")
    logger.error(f"    Found {future_count} future timestamps (should be
0)")

    for error in errors[:5]:
        logger.error(f"        • {error}")
    return False

# Date distribution
logger.info(f"\n📅 Date Distribution (top 10):")
sorted_dates = sorted(date_distribution.items(), key=lambda x: x[1],
reverse=True)[:10]
for date_str, count in sorted_dates:
    pct = (count / total_records) * 100
    logger.info(f"    {date_str}: {count:,} records ({pct:.1f}%)")

# Field statistics
logger.info(f"\n📊 Sensor Field Statistics:")
for field, stats in sorted(fields_summary.items()):
    if stats['count'] > 0:
        avg = stats['sum'] / stats['count']
        logger.info(f"    {field.capitalize()}:")
        logger.info(f"        Range: {stats['min']:.2f} to
{stats['max']:.2f}")
        logger.info(f"        Average: {avg:.2f}")

# Sample records
logger.info(f"\n📄 Sample Records:")
for sample in sample_records[:5]:
    logger.info(f"    Note {sample['note_id']}: {sample['original']} →
{sample['updated']} (temp: {sample['temp']})")

# Errors
if errors:
    logger.warning(f"\n⚠️ Warnings ({len(errors)}):")
    for error in errors[:5]:
        logger.warning(f"        • {error}")

```

```

# Final verdict
logger.info(f"\n{'='*80}")
if future_count == 0 and total_records > 0:
    logger.info("✅ VALIDATION PASSED!")
    logger.info("    • No future timestamps")
    logger.info("    • All records have valid updated_timestamp")
    logger.info("    • Data ready for use")
    logger.info(f"\n📝 Next: Update docker-compose.yml and restart containers")
    return True
else:
    logger.error("❌ VALIDATION FAILED!")
    logger.error("    Please check errors above and re-run transformation")
    return False

def main():
    """Main entry point."""
    script_dir = Path(__file__).parent
    project_root = script_dir.parent.parent
    output_file = project_root / "data" / "data_realtime.csv"

    success = validate_output_file(output_file)

    return 0 if success else 1

if __name__ == "__main__":
    exit(main())

```

docker-compose.yml

```

services:

  redpanda:
    image: redpandadata/redpanda:latest
    container_name: sieis-redpanda
    command:
      - redpanda
      - start
      - --smp 1
      - --memory 512M
      - --overprovisioned
      - --kafka-addr internal://0.0.0.0:9092,external://0.0.0.0:19092
      - --advertise-kafka-addr
    internal://redpanda:9092,external://localhost:19092
    ports:
      - "9092:9092"
      - "19092:19092"

```



```

networks:
  - sieis-network

redpanda-console:
  image: redpandadata/console:latest
  container_name: sieis-console
  environment:
    - KAFKA_BROKERS=redpanda:9092
  ports:
    - "8080:8080"
  depends_on:
    - redpanda
  networks:
    - sieis-network

influxdb3:
  image: influxdb:latest
  container_name: sieis-influxdb3
  ports:
    - "8086:8086"
  environment:
    - DOCKER_INFLUXDB_INIT_MODE=setup
    - DOCKER_INFLUXDB_INIT_USERNAME=admin
    - DOCKER_INFLUXDB_INIT_PASSWORD=password123
    - DOCKER_INFLUXDB_INIT_ORG=sieis
    - DOCKER_INFLUXDB_INIT_BUCKET=sensor_data
    - DOCKER_INFLUXDB_INIT_ADMIN_TOKEN=my-super-secret-token
    - DOCKER_INFLUXDB_INIT_RETENTION=0 # Infinite retention for development
    (use 30d, 90d, 365d for production)
  volumes:
    - influxdb3_data:/var/lib/influxdb2
  networks:
    - sieis-network
  healthcheck:
    test: ["CMD", "curl", "-f", "http://localhost:8086/health"]
    interval: 30s
    timeout: 10s
    retries: 5

minio:
  image: minio/minio:latest
  container_name: sieis-minio
  ports:
    - "9000:9000"
    - "9001:9001"
  environment:
    - MINIO_ROOT_USER=minioadmin
    - MINIO_ROOT_PASSWORD=minioadmin123
  command: server /data --console-address ":9001"
  volumes:

```

```

    - minio_data:/data
networks:
    - sieis-network
healthcheck:
    test: ["CMD", "curl", "-f", "http://localhost:9000/minio/health/live"]
    interval: 30s
    timeout: 10s
    retries: 5

minio-init:
    image: minio/mc:latest
    container_name: sieis-minio-init
    depends_on:
        minio:
            condition: service_healthy
    entrypoint: >
        /bin/sh -c "
        /usr/bin/mc alias set myminio http://minio:9000 minioadmin
minioadmin123;
        /usr/bin/mc mb myminio/sieis-archive --ignore-existing;
        /usr/bin/mc anonymous set download myminio/sieis-archive;
        exit 0;
        "
    networks:
        - sieis-network

simulator:
    build:
        context: .
        dockerfile: Dockerfile.simulator
    container_name: sieis-simulator
    depends_on:
        - redpanda
    environment:
        - KAFKA_BROKER=redpanda:9092
        - KAFKA_TOPIC=sensor_readings
        - SPEED_FACTOR=100
        - DATA_PATH=/app/data/processed/incremental_data.txt
        - MOTE_LOCS_PATH=/app/data/raw/mote_locs.txt
        - FILTER_TODAY_ONLY=true
    volumes:
        - ./data/raw:/app/data/raw:ro
        - ./data/processed:/app/data/processed:ro
    networks:
        - sieis-network

consumer:
    build:
        context: .
        dockerfile: Dockerfile.consumer

```

```
container_name: sieis-consumer
depends_on:
  redpanda:
    condition: service_started
  influxdb3:
    condition: service_healthy
  minio:
    condition: service_healthy
  minio-init:
    condition: service_completed_successfully
environment:
  - KAFKA_BROKER=redpanda:9092
  - KAFKA_TOPIC=sensor_readings
  - INFLUX_URL=http://influxdb3:8086
  - INFLUX_TOKEN=my-super-secret-token
  - INFLUX_ORG=sieis
  - INFLUX_BUCKET=sensor_data
  - MINIO_ENDPOINT=minio:9000
  - MINIO_ACCESS_KEY=minioadmin
  - MINIO_SECRET_KEY=minioadmin123
  - MINIO_BUCKET=sieis-archive
  - MINIO_SECURE=false
networks:
  - sieis-network
```

```
api:
  build:
    context: .
    dockerfile: Dockerfile.api
  container_name: sieis-api
  ports:
    - "8000:8000"
  depends_on:
    influxdb3:
      condition: service_healthy
    minio:
      condition: service_healthy
  environment:
    - INFLUX_URL=http://influxdb3:8086
    - INFLUX_TOKEN=my-super-secret-token
    - INFLUX_ORG=sieis
    - INFLUX_BUCKET=sensor_data
    - MINIO_ENDPOINT=minio:9000
    - MINIO_ACCESS_KEY=minioadmin
    - MINIO_SECRET_KEY=minioadmin123
    - MINIO_BUCKET=sieis-archive
    - MINIO_SECURE=false
    - API_PORT=8000
    - API_HOST=0.0.0.0
```

```

    volumes:
      - ./src:/app/src                                # live-reload source changes
without rebuild
      - ./src/app/ml/models:/app/src/app/ml/models
    networks:
      - sieis-network
    healthcheck:
      test: ["CMD", "python", "-c", "import urllib.request;
urllib.request.urlopen('http://localhost:8000/api/v1/health')"]
      interval: 15s
      timeout: 10s
      retries: 5
      start_period: 60s

```

```

dashboard:
  build:
    context: .
    dockerfile: Dockerfile.dashboard
  container_name: sieis-dashboard
  ports:
    - "8501:8501"
  depends_on:
    - api
    - influxdb3
    - minio
  environment:
    - INFLUX_URL=http://influxdb3:8086
    - INFLUX_TOKEN=my-super-secret-token
    - INFLUX_ORG=sieis
    - INFLUX_BUCKET=sensor_data
    - MINIO_ENDPOINT=minio:9000
    - MINIO_ACCESS_KEY=minioadmin
    - MINIO_SECRET_KEY=minioadmin123
    - MINIO_BUCKET=sieis-archive
    - MINIO_SECURE=false
    - API_PORT=8000
    - API_URL=http://sieis-api:8000

```

```

    volumes:
      - ./src:/app/src:ro                            # live-reload source changes without
rebuild
      - ./data:/app/data:ro
      - ./src/app/ml/models:/app/src/app/ml/models:ro
    networks:
      - sieis-network

```

— SIEIS Scheduler

```

# APScheduler container: runs two daily cron jobs
#   Job A  00:00 UTC  - remap incremental_data.txt + restart simulator
#   Job B  02:00 UTC  - retrain anomaly model + hot-reload FastAPI

```

```

#
# To run both jobs IMMEDIATELY (smoke-test without waiting for midnight):
#   RUN_JOBS_ON_START=true docker-compose up scheduler
#
scheduler:
  build:
    context: .
    dockerfile: Dockerfile.scheduler
  container_name: sieis-scheduler
  depends_on:
    influxdb3:
      condition: service_healthy
    minio:
      condition: service_healthy
    api:
      condition: service_healthy
  environment:
    # — Schedule times (UTC)
    - SCHEDULER_TIMEZONE=UTC
    - JOB_A_HOUR=0          # midnight UTC → remap + restart simulator
    - JOB_A_MINUTE=0
    - JOB_B_HOUR=2          # 2 AM UTC → retrain model + reload API
    - JOB_B_MINUTE=0
    # Set to "true" to fire both jobs immediately on container start
(testing)
    - RUN_JOBS_ON_START=false
    # — Job A config
    - SIMULATOR_CONTAINER=sieis-simulator
    - INCR_DATA_PATH=/app/data/processed/incremental_data.txt
    # — Job B config
    - API_RELOAD_URL=http://sieis-api:8000/api/v1/ml/model/reload
    # — Shared service credentials
    - MINIO_ENDPOINT=minio:9000
    - MINIO_ACCESS_KEY=minioadmin
    - MINIO_SECRET_KEY=minioadmin123
    - MINIO_BUCKET=sieis-archive
    - MINIO_SECURE=false
    - INFLUX_URL=http://influxdb3:8086
    - INFLUX_TOKEN=my-super-secret-token
    - INFLUX_ORG=sieis
    - INFLUX_BUCKET=sensor_data
  volumes:
    # Data volume — scheduler writes remapped incremental_data.txt here
    - ./data/processed:/app/data/processed
    # Model volume — scheduler writes new .pkl here; API reads from same
mount
    - ./src/app/ml/models:/app/src/app/ml/models
    # Docker socket — allows scheduler to restart sieis-simulator
    - /var/run/docker.sock:/var/run/docker.sock
  networks:

```

```
    - sieis-network
restart: unless-stopped
```

```
networks:
  sieis-network:
    driver: bridge
```

```
volumes:
  influxdb3_data:
  minio_data:
```

pytest.ini

```
[pytest]
# Pytest configuration for SIEIS project

# Register custom marks
markers =
    integration: marks tests as integration tests (requires Docker/external
services)

# Test discovery patterns
python_files = test_*.py
python_classes = Test*
python_functions = test_*

# Output options
console_output_style = progress

# Warnings
filterwarnings =
    ignore::DeprecationWarning
```

requirements.txt

```
kafka-python==2.0.2
influxdb-client==1.38.0
influxdb3-python==0.5.0
pandas==2.1.0
scikit-learn==1.3.0
fastapi==0.103.0
uvicorn==0.23.0
streamlit==1.27.0
plotly==5.17.0
python-dotenv==1.0.0
requests==2.31.0
tqdm==4.66.0
minio==7.2.9
pyarrow==15.0.0
polars==0.20.7
joblib>=1.3.0
numpy>=1.24.0
pydantic>=2.0.0
# — Scheduler dependencies
```

```
APScheduler==3.10.4
docker>=6.1.0
```

scripts/check_kafka_messages.py

```
"""Check Kafka message format."""
from kafka import KafkaConsumer
import json
import pprint

consumer = KafkaConsumer(
    'sensor_readings',
    bootstrap_servers='localhost:19092',
    auto_offset_reset='earliest',
    value_deserializer=lambda m: json.loads(m.decode('utf-8')),
    consumer_timeout_ms=2000
)

print("Reading first few messages from Kafka...")
msgs = []
for msg in consumer:
    msgs.append(msg.value)
    if len(msgs) >= 5:
        break

print(f"\nReceived {len(msgs)} messages in 5 seconds")

if msgs:
```

```

        print("\n" + "=" * 80)
        print("SAMPLE MESSAGE FROM KAFKA:")
        print("=" * 80)
        pprint.pprint(msgs[0])
        print("=" * 80)
    else:
        print("No messages received")

```

scripts/check_latest_timestamps.py

```

"""Check the latest updated_timestamp values in realtime_data.txt"""

```

```

import pandas as pd
from pathlib import Path

```

```

# Read data file
data_file = Path(__file__).parent.parent / "data" / "processed" /
"realtime_data.txt"

```

```

print(f"\n📂 Reading: {data_file.name}\n")

```

```

df = pd.read_csv(
    data_file,
    sep=r'\s+',
    header=None,
    names=['date', 'time', 'epoch', 'moteid', 'temperature',
           'humidity', 'light', 'voltage', 'updated_timestamp']
)

```

```

# Parse timestamps
df['updated_timestamp'] = pd.to_datetime(df['updated_timestamp'],
format='ISO8601')

```

```

# Sort by timestamp descending
df_sorted = df.sort_values('updated_timestamp', ascending=False)

```

```

print("="*80)
print("LATEST UPDATED_TIMESTAMP VALUES")
print("="*80)

```

```

print(f"\n⬆️ Last 20 timestamps (most recent first):\n")
print(df_sorted[['moteid',
'updated_timestamp']].head(20).to_string(index=False))

```

```

latest = df_sorted['updated_timestamp'].iloc[0]
oldest = df_sorted['updated_timestamp'].iloc[-1]

```

```

print(f"\n{'='*80}")
print(f"🇩🇪 TIMESTAMP SUMMARY")

```



```

print(f"{' '*80}")
print(f"    Latest timestamp:    {latest}")
print(f"    Oldest timestamp:     {oldest}")
print(f"    Date range:            {oldest.date()} to {latest.date()}")
print(f"    Total records:        {len(df):,}")
print()

```

scripts/clear_storage.py

"""Clear all data from InfluxDB and MinIO for a fresh load.

Usage:

```

python scripts/clear_storage.py
python scripts/clear_storage.py --influx-only
python scripts/clear_storage.py --minio-only
python scripts/clear_storage.py --dry-run

```

WARNING: This is destructive. All sensor data will be deleted.

"""

```

import sys
import os
sys.path.insert(0,
os.path.dirname(os.path.dirname(os.path.abspath(__file__))))

from dotenv import load_dotenv
load_dotenv()

import argparse
import json
from datetime import datetime
from pathlib import Path

from src.app import config

def clear_influxdb(dry_run: bool = False) -> bool:
    """Drop and recreate InfluxDB bucket (instant regardless of record
count)."""
    print("\n🚧 Clearing InfluxDB...")
    try:
        from influxdb_client import InfluxDBClient

        client = InfluxDBClient(
            url=config.INFLUX_URL,
            token=config.INFLUX_TOKEN,
            org=config.INFLUX_ORG,
            timeout=30_000,  # 30s is plenty - drop/create is instant
        )
    
```

```

# Ping first
if not client.ping():
    print(" ❌ Cannot reach InfluxDB")
    return False

# Count existing records before drop
query_api = client.query_api()
flux = f"""
from(bucket: "{config.INFLUX_BUCKET}")
  |> range(start: 2020-01-01T00:00:00Z, stop: 2030-01-01T00:00:00Z)
  |> filter(fn: (r) => r["_measurement"] == "sensor_reading")
  |> filter(fn: (r) => r["_field"] == "temperature")
  |> count()
"""

    try:
        tables = query_api.query(flux, org=config.INFLUX_ORG)
        existing = sum(r.get_value() or 0 for t in tables for r in
t.records)
        print(f"  Records found : {existing:,}")
    except Exception:
        print("  Records found : (count query timed out – proceeding
anyway)")

    if dry_run:
        print("  DRY RUN – no data deleted")
        client.close()
        return True

# — Drop + recreate the bucket (instantaneous, avoids delete-API
timeout) —
buckets_api = client.buckets_api()

    existing_bucket =
buckets_api.find_bucket_by_name(config.INFLUX_BUCKET)
    if existing_bucket is None:
        print(f" ⚠️ Bucket '{config.INFLUX_BUCKET}' not found – creating
it fresh")
    else:
        buckets_api.delete_bucket(existing_bucket)
        print(f"  Dropped bucket : '{config.INFLUX_BUCKET}'")

# Resolve org ID
orgs_api = client.organizations_api()
org_list = orgs_api.find_organizations(org=config.INFLUX_ORG)
if not org_list:
    print(f" ❌ Organisation '{config.INFLUX_ORG}' not found in
InfluxDB")
    client.close()

```

```

        return False
    org_id = org_list[0].id

    buckets_api.create_bucket(bucket_name=config.INFLUX_BUCKET,
org_id=org_id)
    print(f" 🟢 Recreated bucket '{config.INFLUX_BUCKET}' (empty, ready
for fresh load)")

    client.close()
    return True

except Exception as e:
    print(f" 🛑 InfluxDB clear failed: {e}")
    return False

def clear_minio(dry_run: bool = False) -> bool:
    """Delete all objects from MinIO bucket."""
    print(f"\n 🗑 Clearing MinIO...")
    try:
        from minio import Minio

        client = Minio(
            config.MINIO_ENDPOINT,
            access_key=config.MINIO_ACCESS_KEY,
            secret_key=config.MINIO_SECRET_KEY,
            secure=config.MINIO_SECURE,
        )

        if not client.bucket_exists(config.MINIO_BUCKET):
            print(f" ⚠ Bucket '{config.MINIO_BUCKET}' does not exist -
nothing to clear")
            return True

        # List all objects
        objects = list(client.list_objects(config.MINIO_BUCKET,
recursive=True))
        total = len(objects)
        parquet_count = sum(1 for o in objects if
o.object_name.endswith(".parquet"))
        total_mb = sum(o.size or 0 for o in objects) / (1024 * 1024)

        print(f" Objects found : {total:,} ({parquet_count:,} Parquet,
{total_mb:.1f} MB)")

        if dry_run:
            print(" DRY RUN - no objects deleted")
            return True

```

```

    if total == 0:
        print(" ✅ Bucket already empty")
        return True

    # Delete in batches
    deleted = 0
    for obj in objects:
        client.remove_object(config.MINIO_BUCKET, obj.object_name)
        deleted += 1
        if deleted % 100 == 0:
            print(f"   Deleted {deleted}/{total} objects...")

    print(f" ✅ Deleted {deleted:,} objects from bucket
'{config.MINIO_BUCKET}'")
    return True

except Exception as e:
    print(f" ❌ MinIO clear failed: {e}")
    return False

def reset_checkpoint() -> None:
    """Reset the load checkpoint so loader starts from line 0."""
    checkpoint_path = Path("data/.checkpoint_historical.json")
    checkpoint_path.parent.mkdir(parents=True, exist_ok=True)
    with open(checkpoint_path, "w") as f:
        json.dump({
            "last_line": 0,
            "timestamp": datetime.now().isoformat(),
            "stats": {
                "total_read": 0,
                "influx_written": 0,
                "minio_written": 0,
                "skipped": 0,
                "errors": 0,
            },
            "note": "Reset by clear_storage.py",
        }, f, indent=2)
    print(f"\n🔄 Checkpoint reset: {checkpoint_path}")

def main():
    parser = argparse.ArgumentParser(description="Clear InfluxDB and MinIO for
a fresh data load")
    parser.add_argument("--influx-only", action="store_true", help="Clear
InfluxDB only")
    parser.add_argument("--minio-only", action="store_true", help="Clear
MinIO only")
    parser.add_argument("--dry-run", action="store_true", help="Show what

```

```

would be deleted without deleting")
    args = parser.parse_args()

    do_influx = not args.minio_only
    do_minio = not args.influx_only

    print("=" * 55)
    print("SIEIS Storage Clear")
    print(f"Time      : {datetime.now().strftime('%Y-%m-%d %H:%M:%S')}")
    print(f"InfluxDB: {'yes' if do_influx else 'skip'} - {config.INFLUX_URL} /
{config.INFLUX_BUCKET}")
    print(f"MinIO    : {'yes' if do_minio else 'skip'} -
{config.MINIO_ENDPOINT} / {config.MINIO_BUCKET}")
    print(f"Dry run : {args.dry_run}")
    print("=" * 55)

    if not args.dry_run:
        confirm = input("\n⚠ This will DELETE ALL sensor data. Type YES to
continue: ")
        if confirm.strip().upper() != "YES":
            print("Aborted.")
            sys.exit(0)

    influx_ok = clear_influxdb(dry_run=args.dry_run) if do_influx else True
    minio_ok = clear_minio(dry_run=args.dry_run) if do_minio else True

    if not args.dry_run:
        reset_checkpoint()

    print("\n" + "=" * 55)
    print("SUMMARY")
    print("=" * 55)
    print(f" InfluxDB : {'✅ Cleared' if influx_ok else '❌ Failed'}")
    print(f" MinIO    : {'✅ Cleared' if minio_ok else '❌ Failed'}")

    if not args.dry_run and influx_ok and minio_ok:
        print("\n✅ Storage cleared. Ready for fresh load.")
        print(" Next: python scripts/load_historical_data.py")

if __name__ == "__main__":
    main()

```

[scripts/compact_minio_parquet.py](#)

"""

compact_minio_parquet.py - PyArrow-only MinIO Parquet compaction.

Groups the 166k+ small batch files by their Hive partition

(year / month / day / mote_id) and rewrites each partition as a single snappy-compressed Parquet file, then deletes the originals.

Expected outcome

Before : ~166,228 files (~1 KB each)
After : ~3,000 files (~70 KB each)
Speedup: ML training load time drops from ~30 min → ~3 min

Usage

```
# Dry-run (lists what would happen, no changes)
python scripts/compact_minio_parquet.py --dry-run

# Compact everything
python scripts/compact_minio_parquet.py

# Compact only one specific partition key
python scripts/compact_minio_parquet.py --partition
year=2026/month=02/day=10/mote_id=1084

# Compact but keep originals (no delete)
python scripts/compact_minio_parquet.py --keep-originals

# Limit to N partitions (useful for testing)
python scripts/compact_minio_parquet.py --limit 10
"""

import argparse
import io
import logging
import re
import sys
import time
from collections import defaultdict
from pathlib import Path

import pyarrow as pa
import pyarrow.parquet as pq

# — Logging


---


logging.basicConfig(
    level=logging.INFO,
    format="%(asctime)s [%(levelname)s] %(message)s",
    datefmt="%H:%M:%S",
)
log = logging.getLogger(__name__)

# — Path helpers
```

```
PROJECT_ROOT = Path(__file__).parent.parent
sys.path.insert(0, str(PROJECT_ROOT))
```

```
# Partition key pattern: year=YYYY/month=MM/day=DD/mote_id=XXXXX
PARTITION_RE = re.compile(
    r"^(year=\d+/month=\d+/day=\d+/mote_id=\d+)/(.+\.parquet)$"
)
```

```
def _minio_client():
    from minio import Minio
    from src.app import config
```

```
    return (
        Minio(
            config.MINIO_ENDPOINT,
            access_key=config.MINIO_ACCESS_KEY,
            secret_key=config.MINIO_SECRET_KEY,
            secure=config.MINIO_SECURE,
        ),
        config.MINIO_BUCKET,
    )
```

```
# — Core functions
```

```
def list_partitions(client, bucket: str) -> dict[str, list[str]]:
    """Return {partition_key: [object_name, ...]} for all Parquet files."""
    log.info(f"Listing objects in bucket '{bucket}' ...")
    partitions: dict[str, list[str]] = defaultdict(list)
    skipped = 0

    for obj in client.list_objects(bucket, recursive=True):
        name = obj.object_name
        if not name.endswith(".parquet"):
            continue
        m = PARTITION_RE.match(name)
        if m:
            partitions[m.group(1)].append(name)
        else:
            # Could be an already-compacted file at root level — skip.
            skipped += 1

    log.info(
        f"Found {sum(len(v) for v in partitions.values()):,} Parquet files "
        f"across {len(partitions):,} partitions ({skipped} non-partition "
        f"files skipped)"
    )
```

```

return dict(partitions)

def compact_partition(
    client,
    bucket: str,
    partition_key: str,
    object_names: list[str],
    dry_run: bool,
    keep_originals: bool,
) -> dict:
    """
    Read all files in a partition, concatenate with PyArrow,
    write one compacted file back, delete originals.

    Returns a result dict with counts and status.
    """
    result = {
        "partition": partition_key,
        "input_files": len(object_names),
        "input_rows": 0,
        "output_rows": 0,
        "status": "ok",
        "error": None,
    }

    if len(object_names) == 1:
        result["status"] = "skipped_single"
        return result

    # — Read all small files —————
    tables = []
    clock_skew_count = 0

    for obj_name in object_names:
        for attempt in range(3):
            try:
                response = client.get_object(bucket, obj_name)
                data = response.read()
                response.close()
                response.release_conn()
                tables.append(pq.read_table(io.BytesIO(data)))
                break
            except Exception as e:
                err = str(e)
                if "RequestTimeTooSkewed" in err:
                    clock_skew_count += 1
                    break # no point retrying
                if attempt < 2:
                    time.sleep(0.5)

```



```

        else:
            log.warning(f" ⚠ Failed {obj_name}: {e}")
            result["status"] = "partial"

    if not tables:
        result["status"] = "error"
        result["error"] = "No files could be read"
        return result

    if clock_skew_count:
        log.warning(f" ⚠ {clock_skew_count} files skipped due to clock skew
in {partition_key}")
        result["status"] = "partial"

    # — Concatenate —————
    combined = pa.concat_tables(tables, promote_options="default")
    result["input_rows"] = combined.num_rows
    result["output_rows"] = combined.num_rows

    # — Build output object name —————
    # e.g. year=2026/month=02/day=10/mote_id=1084/compacted.parquet
    output_name = f"{partition_key}/compacted.parquet"

    if dry_run:
        log.info(
            f" [DRY-RUN] {partition_key}: "
            f"{len(object_names)} files → 1 file  ({combined.num_rows:,}
rows)"
        )
        result["status"] = "dry_run"
        return result

    # — Write compacted file —————
    buf = io.BytesIO()
    pq.write_table(combined, buf, compression="snappy")
    buf.seek(0)
    data_bytes = buf.getvalue()

    client.put_object(
        bucket,
        output_name,
        io.BytesIO(data_bytes),
        length=len(data_bytes),
        content_type="application/octet-stream",
    )

    # — Delete originals —————
    if not keep_originals:
        for obj_name in object_names:

```

```

        try:
            client.remove_object(bucket, obj_name)
        except Exception as e:
            log.warning(f" ⚠️ Could not delete {obj_name}: {e}")

    return result

# — Main


---



def main():
    parser = argparse.ArgumentParser(
        description="Compact MinIO Parquet files – one file per partition."
    )
    parser.add_argument(
        "--dry-run", action="store_true",
        help="Show what would happen without making any changes",
    )
    parser.add_argument(
        "--keep-originals", action="store_true",
        help="Write compacted files but do NOT delete the originals",
    )
    parser.add_argument(
        "--partition", type=str, default=None,
        help="Compact only this specific partition key (e.g. year=2026/month=02/day=10/mote_id=1084)",
    )
    parser.add_argument(
        "--limit", type=int, default=0,
        help="Stop after compacting this many partitions (0 = all)",
    )
    parser.add_argument(
        "--min-files", type=int, default=2,
        help="Only compact partitions with at least this many files (default: 2)",
    )
    args = parser.parse_args()

    print("=" * 65)
    print("SIEIS – PyArrow MinIO Parquet Compaction")
    mode = "DRY-RUN" if args.dry_run else ("KEEP ORIGINALS" if
args.keep_originals else "COMPACT + DELETE")
    print(f"Mode          : {mode}")
    print(f"Min files   : {args.min_files}")
    print(f"Limit       : {args.limit if args.limit else 'ALL'}")
    if args.partition:
        print(f"Partition   : {args.partition}")
    print("=" * 65)

```

```

client, bucket = _minio_client()

# — Get partition map —————
if args.partition:
    # Verify partition exists
    all_parts = list_partitions(client, bucket)
    if args.partition not in all_parts:
        print(f"✗ Partition '{args.partition}' not found in bucket.")
        sys.exit(1)
    partitions = {args.partition: all_parts[args.partition]}
else:
    partitions = list_partitions(client, bucket)

# Filter by min-files threshold
partitions = {k: v for k, v in partitions.items() if len(v) >=
args.min_files}
print(f"\nPartitions needing compaction : {len(partitions):,}")

if not partitions:
    print("Nothing to compact. Exiting.")
    return

# Apply limit
partition_list = list(partitions.items())
if args.limit:
    partition_list = partition_list[: args.limit]
    print(f"Limited to first {args.limit} partitions.\n")

# — Compact each partition —————
stats = {"ok": 0, "partial": 0, "error": 0, "dry_run": 0,
"skipped_single": 0}
total_input_files = 0
total_input_rows = 0
start = time.time()

for i, (pkey, obj_names) in enumerate(partition_list, 1):
    prefix = f"[{i}/{len(partition_list)}]"
    print(f"{prefix} {pkey} ({len(obj_names)} files)", end=" ... ",
flush=True)

    result = compact_partition(
        client, bucket, pkey, obj_names,
        dry_run=args.dry_run,
        keep_originals=args.keep_originals,
    )

    status = result["status"]
    stats[status] = stats.get(status, 0) + 1

```

```

total_input_files += result["input_files"]
total_input_rows += result["input_rows"]

if status in ("ok", "partial"):
    print(f"✅ {result['input_rows']:,} rows")
elif status == "dry_run":
    print(f"(dry-run)")
elif status == "skipped_single":
    print(f"skipped (already 1 file)")
else:
    print(f"❌ {result['error']}")

elapsed = time.time() - start
print()
print("=" * 65)
print("COMPACTION COMPLETE")
print(f"  Elapsed           : {elapsed:.1f}s  ({elapsed/60:.1f} min)")
print(f"  Partitions         : {len(partition_list):,}")
print(f"  Input files        : {total_input_files:,}")
print(f"  Input rows         : {total_input_rows:,}")
print(f"  ✅ OK              : {stats.get('ok', 0):,}")
print(f"  ⚠️ Partial         : {stats.get('partial', 0):,}")
print(f"  ❌ Errors          : {stats.get('error', 0):,}")
if args.dry_run:
    print(f"\n ⓘ Run without --dry-run to apply changes.")
print("=" * 65)

if __name__ == "__main__":
    main()

```

scripts/load_historical_data.py

"""

Production-level Historical Data Loader

Loads historical sensor data from data/processed/historical_data.txt into both InfluxDB (hot storage) and MinIO (cold storage, Parquet).

Features:

- Batch processing with configurable batch sizes
- Progress tracking with tqdm
- Error handling and retry logic
- Checkpointing for resume capability
- Deduplication
- Validation and statistics
- Memory-efficient streaming for large files

Usage:

```
python scripts/load_historical_data.py
python scripts/load_historical_data.py --batch-size 10000 --resume
python scripts/load_historical_data.py --dry-run
"""

import sys
from pathlib import Path
sys.path.insert(0, str(Path(__file__).parent.parent))

from dotenv import load_dotenv
load_dotenv() # Must run before config import so .env values are available
```

```
import argparse
import json
import logging
from datetime import datetime
from typing import Dict, List
import pandas as pd
from tqdm import tqdm
import pyarrow as pa
import pyarrow.parquet as pq
from minio import Minio
```

```
from influxdb_client import InfluxDBClient, Point, WritePrecision
from influxdb_client.client.write_api import SYNCHRONOUS
```

```
from src.app.config import (
    INFLUX_URL,
    INFLUX_TOKEN,
    INFLUX_ORG,
    INFLUX_BUCKET,
    MINIO_ENDPOINT,
    MINIO_ACCESS_KEY,
    MINIO_SECRET_KEY,
    MINIO_BUCKET,
    MINIO_SECURE,
)
```

```
# Configure logging
logging.basicConfig(
    level=logging.INFO,
    format='%(asctime)s - %(levelname)s - %(message)s'
)
logger = logging.getLogger(__name__)
```

```
class HistoricalDataLoader:
```

```
    """Production-grade historical data loader with dual-write to InfluxDB and
    MinIO."""
```

```

def __init__(
    self,
    data_path: str,
    batch_size: int = 5000,
    checkpoint_file: str = "data/.checkpoint_historical.json",
    dry_run: bool = False,
    minio_only: bool = False,
):
    """Initialize the historical data loader.

    Args:
        data_path: Path to historical_data.txt file
        batch_size: Number of records per batch
        checkpoint_file: Path to checkpoint file for resume capability
        dry_run: If True, only validate data without writing
        minio_only: If True, write to MinIO only – skip InfluxDB entirely
    """
    self.data_path = Path(data_path)
    self.batch_size = batch_size
    self.checkpoint_file = Path(checkpoint_file)
    self.dry_run = dry_run
    self.minio_only = minio_only

    # Statistics
    self.stats = {
        'total_read': 0,
        'influx_written': 0,
        'minio_written': 0,
        'skipped': 0,
        'errors': 0,
        'start_time': None,
        'end_time': None
    }

    # Initialize clients
    if not dry_run:
        if not minio_only:
            self._init_influxdb()
        self._init_minio()

    logger.info(f"Initialized HistoricalDataLoader:")
    logger.info(f"  Data source: {self.data_path}")
    logger.info(f"  Batch size: {self.batch_size}")
    logger.info(f"  Dry run: {self.dry_run}")

def _init_influxdb(self):
    """Initialize InfluxDB client."""
    logger.info(f"Connecting to InfluxDB at {INFLUX_URL}")
    self.influx_client = InfluxDBClient(

```

```

        url=INFLUX_URL,
        token=INFLUX_TOKEN,
        org=INFLUX_ORG
    )
    self.influx_write_api =
self.influx_client.write_api(write_options=SYNCHRONOUS)

    # Verify connection
    health = self.influx_client.health()
    logger.info(f"InfluxDB health: {health.status}")

def _init_minio(self):
    """Initialize MinIO client."""
    logger.info(f"Connecting to MinIO at {MINIO_ENDPOINT}")
    self.minio_client = Minio(
        endpoint=MINIO_ENDPOINT,
        access_key=MINIO_ACCESS_KEY,
        secret_key=MINIO_SECRET_KEY,
        secure=MINIO_SECURE
    )

    # Ensure bucket exists
    if not self.minio_client.bucket_exists(MINIO_BUCKET):
        self.minio_client.make_bucket(MINIO_BUCKET)
        logger.info(f"Created MinIO bucket: {MINIO_BUCKET}")
    else:
        logger.info(f"MinIO bucket exists: {MINIO_BUCKET}")

def load_checkpoint(self) -> int:
    """Load checkpoint to resume from last position.

    Returns:
        Line number to start from (0 if no checkpoint)
    """
    if not self.checkpoint_file.exists():
        return 0

    try:
        with open(self.checkpoint_file, 'r') as f:
            checkpoint = json.load(f)
            logger.info(f"Resuming from checkpoint: line
{checkpoint['last_line']:,}")
            return checkpoint['last_line']
    except Exception as e:
        logger.warning(f"Failed to load checkpoint: {e}. Starting from
beginning.")
        return 0

def save_checkpoint(self, line_number: int):
    """Save checkpoint for resume capability."""

```

```

# Create serializable stats
serializable_stats = self.stats.copy()
if isinstance(serializable_stats.get('start_time'), datetime):
    serializable_stats['start_time'] =
serializable_stats['start_time'].isoformat()
if isinstance(serializable_stats.get('end_time'), datetime):
    serializable_stats['end_time'] =
serializable_stats['end_time'].isoformat()

```

```

checkpoint = {
    'last_line': line_number,
    'timestamp': datetime.now().isoformat(),
    'stats': serializable_stats
}
self.checkpoint_file.parent.mkdir(parents=True, exist_ok=True)
with open(self.checkpoint_file, 'w') as f:
    json.dump(checkpoint, f, indent=2)

```

```

def parse_line(self, line: str) -> Dict:
    """Parse a single line from historical_data.txt.

```

Actual Intel Lab format: date time epoch mote_id temperature humidity
light voltage [updated_timestamp]

Note: Some lines have missing sensor values or no updated_timestamp.

Examples:

```

- No sensor: 2004-02-28 00:58:15.315133 1 2026-01-19T19:03:19.891610
(4 parts)
- No updated_ts: 2004-02-28 00:58:46.002832 65333 28 19.0 19.7336
37.0933 71.76 (8 parts)
- Full: 2004-02-28 00:58:46.002832 65333 28 19.0 19.7336 37.0933 71.76
2.69964 2026-01-19T19:03:50.579309 (9 parts→ BUG was using parts[2]=epoch as
mote_id)

```

Returns:

Dictionary with parsed fields, or None if parse fails

"""

try:

```

    parts = line.strip().split()

```

Minimum required: date, time, at least mote_id +
updated_timestamp

```

    if len(parts) < 4:
        return None

```

```

    date = parts[0]    # 2004-02-28
    time = parts[1]    # 00:58:15.315133

```

```

    if len(parts) == 4:

```

Sparse line: date time mote_id updated_timestamp (no epoch,


```

no sensors)

        return {
            'mote_id': int(parts[2]),
            'timestamp': f"{date} {time}",
            'updated_timestamp': parts[3],
            'temperature': None,
            'humidity': None,
            'light': None,
            'voltage': None,
        }
    elif len(parts) == 9:
        # Full record: date time epoch mote_id temp humidity light
        voltage updated_timestamp
        # FIX: parts[2]=epoch (skip), parts[3]=mote_id,
        parts[4..7]=sensors, parts[8]=updated_ts
        return {
            'mote_id': int(float(parts[3])),
            'timestamp': f"{date} {time}",
            'updated_timestamp': parts[8],
            'temperature': float(parts[4]) if parts[4] != '?' else
None,
            'humidity': float(parts[5]) if parts[5] != '?' else
None,
            'light': float(parts[6]) if parts[6] != '?' else
None,
            'voltage': float(parts[7]) if parts[7] != '?' else
None,
        }
    elif len(parts) == 8:
        # No updated_timestamp: date time epoch mote_id temp humidity
        light voltage
        # These lack a remapped timestamp – skip to avoid stale 2004
        data in stores
        logger.debug(f"Skipping 8-part line (no updated_timestamp):
{line.strip()[:100]}")
        return None
    else:
        # Unexpected format – log and skip
        logger.debug(f"Unexpected format ({len(parts)} parts):
{line.strip()[:100]}")
        return None

except (ValueError, IndexError) as e:
    logger.debug(f"Failed to parse line: {line.strip()[:100]} – Error:
{e}")
    return None

def batch_to_influx_points(self, batch: List[Dict]) -> List[Point]:
    """Convert batch of records to InfluxDB Points.

```

```

Args:
    batch: List of record dictionaries

Returns:
    List of InfluxDB Point objects
"""
points = []
for record in batch:
    try:
        # Create point
        point = Point("sensor_reading").tag("mote_id",
str(record['mote_id']))

        # Add fields (only non-None values)
        field_count = 0
        for field in ('temperature', 'humidity', 'light', 'voltage'):
            if record.get(field) is not None:
                point = point.field(field, float(record[field]))
                field_count += 1

        # Skip records with no sensor values
        if field_count == 0:
            self.stats['skipped'] += 1
            continue

        # Use updated_timestamp (mapped to current year)
        ts_str = record.get('updated_timestamp')
        if ts_str:
            ts = datetime.fromisoformat(ts_str)
            point = point.time(ts, WritePrecision.NS)
            points.append(point)

    except Exception as e:
        logger.debug(f"Failed to create point: {e}")
        self.stats['skipped'] += 1

return points

def write_batch_to_influx(self, points: List[Point]) -> bool:
    """Write batch to InfluxDB.

Args:
    points: List of InfluxDB Points

Returns:
    True if successful, False otherwise
"""
    if not points:
        return True

```

```

try:
    self.influx_write_api.write(
        bucket=INFLUX_BUCKET,
        org=INFLUX_ORG,
        record=points
    )
    self.stats['influx_written'] += len(points)
    return True
except Exception as e:
    logger.error(f"InfluxDB write failed: {e}")
    self.stats['errors'] += 1
    return False

def write_batch_to_minio(self, batch: List[Dict]) -> bool:
    """Write batch to MinIO as Parquet.

    Groups by date and mote_id for efficient partitioning.

    Args:
        batch: List of record dictionaries

    Returns:
        True if successful, False otherwise
    """
    if not batch:
        return True

    try:
        # Convert to DataFrame
        df = pd.DataFrame(batch)

        # Parse updated_timestamp for partitioning
        df['updated_timestamp'] = pd.to_datetime(df['updated_timestamp'])
        df['year'] = df['updated_timestamp'].dt.year
        df['month'] = df['updated_timestamp'].dt.month
        df['day'] = df['updated_timestamp'].dt.day

        # Group by date and mote_id
        for (year, month, day, mote_id), group in df.groupby(['year',
            'month', 'day', 'mote_id']):
            # Create partition path
            partition_path =
f"year={year}/month={month:02d}/day={day:02d}/mote_id={mote_id}"

            # Generate filename with timestamp
            timestamp = datetime.now().strftime("%Y%m%d_%H%M%S_%f")
            filename = f"{partition_path}/batch_{timestamp}.parquet"

            # Convert to Parquet
            table = pa.Table.from_pandas(group[

```

```

        'mote_id', 'timestamp', 'updated_timestamp',
        'temperature', 'humidity', 'light', 'voltage'
    ]])

    # Write to bytes buffer
    import io
    buffer = io.BytesIO()
    pq.write_table(table, buffer, compression='snappy')
    buffer.seek(0)

    # Upload to MinIO
    self.minio_client.put_object(
        bucket_name=MINIO_BUCKET,
        object_name=filename,
        data=buffer,
        length=buffer.getbuffer().nbytes,
        content_type='application/octet-stream'
    )

    self.stats['minio_written'] += len(batch)
    return True

except Exception as e:
    logger.error(f"MinIO write failed: {e}")
    self.stats['errors'] += 1
    return False

def load(self, resume: bool = False):
    """Load historical data with dual-write to InfluxDB and MinIO.

    Args:
        resume: If True, resume from last checkpoint
    """
    start_line = self.load_checkpoint() if resume else 0
    self.stats['start_time'] = datetime.now()

    logger.info("="*80)
    logger.info("HISTORICAL DATA LOADING")
    logger.info("="*80)
    logger.info(f"Source: {self.data_path}")
    if self.minio_only:
        logger.info(f"Target: MinIO ONLY ({MINIO_BUCKET}) - InfluxDB
skipped")
    else:
        logger.info(f"Target: InfluxDB ({INFLUX_BUCKET}) + MinIO
({MINIO_BUCKET})")
        logger.info(f"Batch size: {self.batch_size}")
        logger.info(f"Resume from line: {start_line:,}")
        logger.info("="*80)

```

```

# Count total lines for progress bar
logger.info("Counting total lines...")
with open(self.data_path, 'r') as f:
    total_lines = sum(1 for _ in f)
logger.info(f"Total lines: {total_lines:,}")

# Process file in batches
batch = []
current_line = 0

with open(self.data_path, 'r') as f:
    # Create progress bar
    pbar = tqdm(total=total_lines, desc="Loading data", unit=" lines")

    for line_num, line in enumerate(f, 1):
        # Skip to checkpoint
        if line_num <= start_line:
            pbar.update(1)
            continue

        # Parse line
        record = self.parse_line(line)
        if record:
            batch.append(record)
            self.stats['total_read'] += 1
        else:
            self.stats['skipped'] += 1

        current_line = line_num
        pbar.update(1)

    # Write batch when full
    if len(batch) >= self.batch_size:
        if not self.dry_run:
            # Write to InfluxDB (skipped in minio_only mode)
            if not self.minio_only:
                points = self.batch_to_influx_points(batch)
                self.write_batch_to_influx(points)

            # Write to MinIO (always)
            self.write_batch_to_minio(batch)

        # Save checkpoint every 10 batches
        if (line_num // self.batch_size) % 10 == 0:
            self.save_checkpoint(current_line)

        batch = []

# Write remaining records
if batch and not self.dry_run:

```

```

        if not self.minio_only:
            points = self.batch_to_influx_points(batch)
            self.write_batch_to_influx(points)
            self.write_batch_to_minio(batch)

    pbar.close()

    # Final checkpoint
    self.save_checkpoint(current_line)
    self.stats['end_time'] = datetime.now()

    # Print summary
    self.print_summary()

def print_summary(self):
    """Print loading summary."""
    duration = (self.stats['end_time'] -
self.stats['start_time']).total_seconds()

    print("\n" + "="*80)
    print("LOADING SUMMARY")
    print("="*80)
    print(f"Duration: {duration:.2f} seconds")
    print(f"Total read: {self.stats['total_read']:,} records")
    print(f"InfluxDB written: {self.stats['influx_written']:,} points")
    print(f"MinIO written: {self.stats['minio_written']:,} records")
    print(f"Skipped: {self.stats['skipped']:,} records")
    print(f"Errors: {self.stats['errors']:,}")
    print(f"Throughput: {self.stats['total_read'] / duration:.0f}
records/sec")
    print("="*80)

    if self.dry_run:
        print("DRY RUN: No data was written")
    else:
        print(f"✅ Data successfully loaded to InfluxDB and MinIO")
    print()

def close(self):
    """Close all connections."""
    if not self.dry_run:
        if hasattr(self, 'influx_client'):
            self.influx_client.close()
        logger.info("Connections closed")

def main():
    """Main entry point."""
    parser = argparse.ArgumentParser(description="Load historical data to

```

```

InfluxDB and MinIO")
    parser.add_argument(
        '--data-file',
        default='data/processed/historical_data.txt',
        help='Path to historical data file (default:
data/processed/historical_data.txt)'
    )
    parser.add_argument(
        '--batch-size',
        type=int,
        default=5000,
        help='Batch size for processing (default: 5000)'
    )
    parser.add_argument(
        '--resume',
        action='store_true',
        help='Resume from last checkpoint'
    )
    parser.add_argument(
        '--dry-run',
        action='store_true',
        help='Validate data without writing'
    )
    parser.add_argument(
        '--minio-only',
        action='store_true',
        help='Write to MinIO only - skip InfluxDB (use for historical batch
load)'
    )

    args = parser.parse_args()

    # Validate data file exists
    data_path = Path(args.data_file)
    if not data_path.exists():
        logger.error(f"Data file not found: {data_path}")
        return 1

    # Initialize loader
    loader = HistoricalDataLoader(
        data_path=str(data_path),
        batch_size=args.batch_size,
        dry_run=args.dry_run,
        minio_only=args.minio_only,
    )

    try:
        # Load data
        loader.load(resume=args.resume)
        return 0

```

```

except KeyboardInterrupt:
    logger.warning("\n\nInterrupted by user")
    logger.info("Progress saved in checkpoint. Run with --resume to
continue.")
    return 1

except Exception as e:
    logger.exception(f"Fatal error: {e}")
    return 1

finally:
    loader.close()

if __name__ == "__main__":
    sys.exit(main())

```

scripts/load_historical_data_new.py

"""

load_historical_data_new.py – Compact-First Historical Data Loader

Reuses the proven parse_line() from load_historical_data.py.

Writes ONE Parquet file per day (all motes combined) to MinIO.

Result:

Before : 166,228 files (~1 KB each)

After : ~65 files (~500 KB each, 1 per day)

ML load: 30 min → seconds

Usage:

python scripts/load_historical_data_new.py

python scripts/load_historical_data_new.py --dry-run

python scripts/load_historical_data_new.py --clear-first

"""

import sys

import os

import io

import logging

import argparse

import json

from collections import defaultdict

from datetime import datetime

from pathlib import Path

sys.path.insert(0, str(Path(__file__).parent.parent))

from dotenv import load_dotenv


```

load_dotenv()

import pandas as pd
import pyarrow as pa
import pyarrow.parquet as pq
from minio import Minio
from tqdm import tqdm

from src.app.config import (
    MINIO_ENDPOINT,
    MINIO_ACCESS_KEY,
    MINIO_SECRET_KEY,
    MINIO_BUCKET,
    MINIO_SECURE,
)

logging.basicConfig(
    level=logging.INFO,
    format="%(asctime)s [%(levelname)s] %(message)s"
)
logger = logging.getLogger(__name__)

PROJECT_ROOT = Path(__file__).resolve().parent.parent
DATA_FILE = PROJECT_ROOT / "data" / "processed" / "historical_data.txt"

# PyArrow schema
SCHEMA = pa.schema([
    pa.field("mote_id", pa.int32()),
    pa.field("temperature", pa.float32()),
    pa.field("humidity", pa.float32()),
    pa.field("light", pa.float32()),
    pa.field("voltage", pa.float32()),
    pa.field("timestamp", pa.string()),
    pa.field("updated_timestamp", pa.string()),
])

# — Parser


---


def parse_line(line: str) -> dict:
    """
    Parse a single line from historical_data.txt.

    Actual Intel Lab format:
        4 parts : date time mote_id updated_timestamp          (no epoch, no
sensors)
        9 parts : date time epoch mote_id temp hum light volt updated_timestamp

    NOTE: For 9-part lines parts[2] is the epoch sequence number (NOT
mote_id).

```

```

        mote_id is at parts[3], sensors at parts[4-7], updated_ts at
parts[8].
        """
        try:
            parts = line.strip().split()

            if len(parts) < 4:
                return None

            date = parts[0]    # 2004-02-28
            time = parts[1]    # 00:58:46.002832

            if len(parts) == 4:
                # Sparse line: date time mote_id updated_timestamp (no epoch
column)
                return {
                    "mote_id":          int(parts[2]),
                    "timestamp":        f"{date} {time}",
                    "updated_timestamp": parts[3],
                    "temperature":      None,
                    "humidity":         None,
                    "light":            None,
                    "voltage":          None,
                }

            elif len(parts) >= 9:
                # Full record: date time epoch mote_id temp hum light volt
updated_ts
                # parts[2] = epoch (skip), parts[3] = mote_id, parts[4-7] =
sensors
                return {
                    "mote_id":          int(float(parts[3])),
                    "timestamp":        f"{date} {time}",
                    "updated_timestamp": parts[8],
                    "temperature":      float(parts[4]) if parts[4] != "?" else
None,
                    "humidity":         float(parts[5]) if parts[5] != "?" else
None,
                    "light":            float(parts[6]) if parts[6] != "?" else
None,
                    "voltage":          float(parts[7]) if parts[7] != "?" else
None,
                }

            else:
                return None

        except (ValueError, IndexError):
            return None

```

— MinIO helpers

```
def get_minio_client() -> Minio:
    return Minio(
        MINIO_ENDPOINT,
        access_key=MINIO_ACCESS_KEY,
        secret_key=MINIO_SECRET_KEY,
        secure=MINIO_SECURE,
    )

def clear_bucket(client: Minio, bucket: str):
    objects = list(client.list_objects(bucket, recursive=True))
    if not objects:
        logger.info("Bucket already empty")
        return
    logger.info(f"Deleting {len(objects):,} existing objects...")
    for obj in objects:
        client.remove_object(bucket, obj.object_name)
    logger.info(f"✅ Deleted {len(objects):,} objects")

def upload_day_parquet(
    client: Minio,
    bucket: str,
    year: int,
    month: int,
    day: int,
    records: list,
) -> int:
    """Compact all records for one day → upload as single Parquet file."""
    df = pd.DataFrame(records)

    # Drop rows with no sensor values
    df = df.dropna(subset=["temperature", "humidity", "light", "voltage"])
    if df.empty:
        return 0

    # Cast types
    df["mote_id"] = df["mote_id"].astype("int32")
    df["temperature"] = df["temperature"].astype("float32")
    df["humidity"] = df["humidity"].astype("float32")
    df["light"] = df["light"].astype("float32")
    df["voltage"] = df["voltage"].astype("float32")
    df["timestamp"] = df["timestamp"].astype(str)
    df["updated_timestamp"] = df["updated_timestamp"].astype(str)

    table = pa.Table.from_pandas(
```

```

df[["mote_id", "temperature", "humidity", "light", "voltage", "timestamp", "updated_
timestamp"]],
    schema=SCHEMA,
    preserve_index=False,
)

buf = io.BytesIO()
pq.write_table(table, buf, compression="snappy", use_dictionary=True)
buf.seek(0)
size = buf.getbuffer().nbytes

object_name = f"year={year}/month={month:02d}/day={day:02d}/data.parquet"
client.put_object(
    bucket,
    object_name,
    buf,
    length=size,
    content_type="application/octet-stream",
)
return size

```

— Main

```

def main():
    parser = argparse.ArgumentParser(
        description="Load historical data into MinIO - compact-first (1 file
per day)"
    )
    parser.add_argument("--data-file", default=str(DATA_FILE))
    parser.add_argument("--dry-run", action="store_true", help="Preview
only, no uploads")
    parser.add_argument("--clear-first", action="store_true", help="Delete
all existing MinIO objects first")
    args = parser.parse_args()

    print("=" * 60)
    print("SIEIS - Historical Data Loader (Compact-First)")
    print(f"Source : {args.data_file}")
    print(f"Bucket : {MINIO_BUCKET}")
    print(f"Mode   : {'DRY-RUN' if args.dry_run else 'LIVE'}")
    print("=" * 60)

```

— Connect

```

print("\n[1/4] Connecting to MinIO...")
client = get_minio_client()
if not client.bucket_exists(MINIO_BUCKET):
    client.make_bucket(MINIO_BUCKET)
    logger.info(f"Created bucket: {MINIO_BUCKET}")

```

```

else:
    logger.info(f"Bucket exists: {MINIO_BUCKET}")

if args.clear_first and not args.dry_run:
    print("\n[!] Clearing existing bucket contents...")
    clear_bucket(client, MINIO_BUCKET)

# — Parse —————
print("\n[2/4] Parsing historical data...")
data_path = Path(args.data_file)
if not data_path.exists():
    print(f"❌ File not found: {data_path}")
    sys.exit(1)

# Count lines for progress bar
with open(data_path, "r") as f:
    total_lines = sum(1 for _ in f)
logger.info(f"Total lines: {total_lines:,}")

# Group records by (year, month, day) using updated_timestamp
day_buckets = defaultdict(list) # key: (year, month, day)
parsed = skipped = 0

with open(data_path, "r") as f:
    for line in tqdm(f, total=total_lines, desc="Parsing", unit=" lines"):
        record = parse_line(line)
        if record is None or record["updated_timestamp"] is None:
            skipped += 1
            continue
        # Only keep records with sensor values
        if record["temperature"] is None:
            skipped += 1
            continue
        try:
            ts = datetime.fromisoformat(record["updated_timestamp"])
            key = (ts.year, ts.month, ts.day)
            day_buckets[key].append(record)
            parsed += 1
        except ValueError:
            skipped += 1

print(f"\n✅ Parsed : {parsed:,} rows with sensor data")
print(f"🔊 Skipped : {skipped:,} rows (no sensor values or parse error)")
print(f"📅 Days : {len(day_buckets)} unique day partitions")

if parsed == 0:
    print(f"❌ No data parsed.")
    sys.exit(1)

```

```

# — Compact + Upload —————
print("\n[3/4] Compacting and uploading (1 file per day)...")

stats = {"files": 0, "rows": 0, "bytes": 0}

for (year, month, day), records in tqdm(
    sorted(day_buckets.items()),
    desc="Uploading",
    unit=" days",
):
    if args.dry_run:
        motes = len({r["mote_id"] for r in records})
        logger.info(
            f" [DRY-RUN]
year={year}/month={month:02d}/day={day:02d}/data.parquet "
            f"— {len(records):,} rows, {motes} motes"
        )
        stats["files"] += 1
        stats["rows"] += len(records)
        continue

    try:
        size = upload_day_parquet(client, MINIO_BUCKET, year, month, day,
records)
        if size > 0:
            stats["files"] += 1
            stats["rows"] += len(records)
            stats["bytes"] += size
    except Exception as e:
        logger.error(f"Failed year={year}/month={month:02d}/day={day:02d}:
{e}")

# — Summary —————
print("\n[4/4] Summary")
print("=" * 60)
print(f" Rows parsed      : {parsed:,}")
print(f" Day partitions    : {stats['files']}")
print(f" Rows uploaded      : {stats['rows']: ,}")
print(f" Data uploaded      : {stats['bytes']/1024/1024:.2f} MB")
print("=" * 60)

if args.dry_run:
    print("\n🚧 DRY-RUN — no changes made. Remove --dry-run to execute.")
else:
    print("\n✅ Done! Next step:")
    print("    python scripts/train_model.py --source minio")

```

```
if __name__ == "__main__":
    main()
```

scripts/ml/diagnose_model.py

```
import pickle, json, pandas as pd
from pathlib import Path
from datetime import datetime

reg = json.load(open('src/app/ml/models/model_registry.json'))
print('Full registry:')
print(json.dumps(reg, indent=2))

# Find the model path from whatever structure exists
model_path = None

# Structure 1: reg['models'][version]['path']
if 'models' in reg:
    models = reg['models']
    latest_key = reg.get('latest')
    print('Latest key:', latest_key)
    if latest_key and latest_key in models:
        model_path = models[latest_key]['path']
    else:
        # Get last model
        last = sorted(models.keys())[-1]
        model_path = models[last]['path']
        print('Using last model:', last)

# Structure 2: reg[version]['path']
elif any('path' in v for v in reg.values() if isinstance(v, dict)):
    for k, v in sorted(reg.items()):
        if isinstance(v, dict) and 'path' in v:
            model_path = v['path']
            print('Found model path:', model_path)

print('Model path:', model_path)

if not model_path:
    print('ERROR: Could not find model path in registry')
    exit(1)

# Try both relative and absolute
p = Path(model_path)
if not p.exists():
    p = Path('src/app/ml/models') / p.name
    print('Trying relative path:', p)

if not p.exists():
```

```

print('ERROR: Model file not found at', p)
# List all pickle files
print('Available pickle files:')
for f in Path('src/app/ml/models').glob('*.pkl'):
    print(' ', f)
exit(1)

print('Loading model from:', p)
with open(p, 'rb') as f:
    artifact = pickle.load(f)

print('Artifact type:', type(artifact))
if isinstance(artifact, dict):
    print('Artifact keys:', list(artifact.keys()))
    pipeline = artifact['pipeline']
    print('Metrics:', json.dumps(artifact.get('metrics'), indent=2))
else:
    pipeline = artifact

now = datetime.utcnow()
test_data = pd.DataFrame([
    {'temperature': 22.5, 'humidity': 55.0, 'light': 300.0, 'voltage': 2.9,
     'hour': now.hour, 'day_of_week': now.weekday()},
    {'temperature': 20.0, 'humidity': 60.0, 'light': 200.0, 'voltage': 3.0,
     'hour': now.hour, 'day_of_week': now.weekday()},
    {'temperature': 25.0, 'humidity': 50.0, 'light': 400.0, 'voltage': 2.8,
     'hour': now.hour, 'day_of_week': now.weekday()},
    {'temperature': 0.0, 'humidity': 0.0, 'light': 0.0, 'voltage': 0.0,
     'hour': now.hour, 'day_of_week': now.weekday()},
    {'temperature': 95.0, 'humidity': 2.0, 'light': 0.0, 'voltage': 0.1,
     'hour': now.hour, 'day_of_week': now.weekday()},
])

preds = pipeline.predict(test_data)
scores = pipeline.decision_function(test_data)

print('\n--- Predictions ---')
labels = ['normal_22.5', 'normal_20.0', 'normal_25.0', 'ZEROS', 'EXTREME']
for label, pred, score in zip(labels, preds, scores):
    flag = 'ANOMALY ❌' if pred == -1 else 'normal ✅'
    print(f' {label:15s} pred={pred:d} score={score:+.4f} -> {flag}')

```


scripts/ml/diagnose_registry.py

```
import pickle, json, pandas as pd
from pathlib import Path
from datetime import datetime

reg = json.load(open('src/app/ml/models/model_registry.json'))
print('Full registry:')
print(json.dumps(reg, indent=2))
```

scripts/remap_timestamps.py

"""Remap updated_timestamp fields in historical and incremental data files.

Strategy

Historical 80%

- Find the max updated_timestamp in the file → call it hist_last_ts
- Compute offset = yesterday - hist_last_ts.date()
- Add that offset to EVERY updated_timestamp in the file
- Result: last historical record lands on yesterday

Incremental 20%

- Collect all unique ORIGINAL dates (parts[0], e.g. 2004-03-21)
- Sort them → map unique_dates[0] → today, unique_dates[1] → today+1, ...
- For each line: keep the time-of-day, swap only the date
- Result: each original calendar day becomes one real calendar day from today

Analogy: Historical is like a printed book – shift all its page-dates by the same number of days so the last page says "yesterday". Incremental is like a live diary – each chapter (original date) gets a new publication date starting from today.

Usage

```
python scripts/remap_timestamps.py          # remap both files
python scripts/remap_timestamps.py --hist-only
python scripts/remap_timestamps.py --incr-only
python scripts/remap_timestamps.py --dry-run  # preview, no writes
"""
```

```
import sys
import os
sys.path.insert(0,
os.path.dirname(os.path.dirname(os.path.abspath(__file__))))

from dotenv import load_dotenv
load_dotenv()
```

```

import argparse
import shutil
from datetime import date, datetime, timedelta
from pathlib import Path

# — helpers


---



def _last_field(parts: list) -> str | None:
    """Return updated_timestamp (last column) if the line has valid data."""
    if len(parts) in (4, 9):
        return parts[-1]
    return None

def _replace_last_field(line: str, new_value: str) -> str:
    """Swap the last whitespace-separated token in a line, keeping all
    others."""
    parts = line.rstrip("\r\n").rsplit(None, 1)
    if len(parts) == 2:
        # Preserve original spacing prefix
        prefix = parts[0]
        return prefix + " " + new_value + "\n"
    return line # malformed — leave untouched

def _parse_ts(ts_str: str) -> datetime | None:
    """Parse ISO timestamp string, tolerating Z suffix."""
    try:
        return datetime.fromisoformat(ts_str.replace("Z",
"+00:00").split("+")[0])
    except Exception:
        return None

# — scan passes


---



def _find_hist_max_date(path: Path) -> date | None:
    """Single-pass scan to find the latest updated_timestamp date in the
    file."""
    max_date = None
    with open(path, "r", encoding="utf-8", errors="replace") as f:
        for line in f:
            parts = line.split()
            ts_str = _last_field(parts)
            if not ts_str:
                continue
            ts = _parse_ts(ts_str)

```

```

        if ts and (max_date is None or ts.date() > max_date):
            max_date = ts.date()
    return max_date

def _find_incr_orig_dates(path: Path) -> list:
    """Collect unique original dates (parts[0]) from incremental file,
    sorted."""
    dates = set()
    with open(path, "r", encoding="utf-8", errors="replace") as f:
        for line in f:
            parts = line.split()
            if len(parts) >= 4:
                try:
                    dates.add(date.fromisoformat(parts[0]))
                except ValueError:
                    pass
    return sorted(dates)

# — remap functions


---



def remap_historical(path: Path, dry_run: bool = False) -> dict:
    """Shift all updated_timestamps in historical file so max date =
    yesterday."""
    print(f"\n[Historical] Scanning {path.name} for max updated_timestamp...")
    hist_max = _find_hist_max_date(path)
    if hist_max is None:
        print(" ❌ Could not find any valid timestamps.")
        return {}

    yesterday = date.today() - timedelta(days=1)
    offset = timedelta(days=(yesterday - hist_max).days)

    print(f" Max date found : {hist_max}")
    print(f" Yesterday      : {yesterday}")
    print(f" Offset applied : {'+' if offset.days >= 0 else ''}{offset.days}
    days")

    if dry_run:
        print(" DRY RUN — no file written.")
        return {"max_before": str(hist_max), "max_after": str(yesterday),
        "offset_days": offset.days}

    # Backup
    backup = path.with_suffix(".txt.bak")
    shutil.copy2(path, backup)
    print(f" Backup written : {backup.name}")

```

```

# Stream-rewrite
tmp = path.with_suffix(".txt.tmp")
changed = skipped = 0
with open(path, "r", encoding="utf-8", errors="replace") as src, \
    open(tmp, "w", encoding="utf-8") as dst:
    for line in src:
        parts = line.split()
        ts_str = _last_field(parts)
        if ts_str:
            ts = _parse_ts(ts_str)
            if ts:
                new_ts = (ts + offset).isoformat(timespec="microseconds")
                # isoformat uses +00:00 suffix if tz-aware; strip it for
consistency
                new_ts = new_ts.split("+")[0]
                dst.write(_replace_last_field(line, new_ts))
                changed += 1
                continue
            dst.write(line)
            skipped += 1

    tmp.replace(path)
    print(f"  Lines remapped : {changed:}, (skipped {skipped:},
unparseable)")
    print(f"  ✅ {path.name} updated.")
    return {"max_before": str(hist_max), "max_after": str(yesterday),
            "offset_days": offset.days, "changed": changed}

def remap_incremental(path: Path, dry_run: bool = False) -> dict:
    """Remap each unique original date in incremental file → today + N
days."""
    print(f"\n[Incremental] Scanning {path.name} for unique original dates...")
    orig_dates = _find_incr_orig_dates(path)
    if not orig_dates:
        print(f"  ❌ No valid dates found.")
        return {}

    today = date.today()
    date_map = {d: today + timedelta(days=i) for i, d in
enumerate(orig_dates)}

    print(f"  Unique original dates : {len(orig_dates)}")
    print(f"  Mapping preview:")
    for orig, new in list(date_map.items())[:3]:
        print(f"    {orig} → {new}")
    if len(date_map) > 3:
        last_orig, last_new = list(date_map.items())[-1]

```

```

print(f"    ...")
print(f"    {last_orig} → {last_new} (last)")

if dry_run:
    print(" DRY RUN – no file written.")
    return {
        "unique_dates": len(orig_dates),
        "first_day": str(today),
        "last_day": str(today + timedelta(days=len(orig_dates) - 1)),
    }

# Backup
backup = path.with_suffix(".txt.bak")
shutil.copy2(path, backup)
print(f" Backup written : {backup.name}")

# Stream-rewrite
tmp      = path.with_suffix(".txt.tmp")
changed = skipped = 0
with open(path, "r", encoding="utf-8", errors="replace") as src, \
    open(tmp, "w", encoding="utf-8") as dst:
    for line in src:
        parts = line.split()
        ts_str = _last_field(parts)
        if ts_str and len(parts) >= 4:
            try:
                orig_date = date.fromisoformat(parts[0])
                new_date  = date_map.get(orig_date)
                ts        = _parse_ts(ts_str)
                if new_date and ts:
                    # Keep the exact time-of-day, swap only the date
                    new_ts = datetime.combine(new_date,
ts.time()).isoformat(timespec="microseconds")
                    dst.write(_replace_last_field(line, new_ts))
                    changed += 1
                    continue
            except (ValueError, KeyError):
                pass
            dst.write(line)
            skipped += 1

tmp.replace(path)
print(f" Lines remapped : {changed:}, (skipped {skipped:},
unparseable)")
print(f"  ✓ {path.name} updated.")
return {
    "unique_dates": len(orig_dates),
    "first_day": str(today),
    "last_day": str(today + timedelta(days=len(orig_dates) - 1)),

```

```

        "changed": changed,
    }

# — main


---



def main():
    parser = argparse.ArgumentParser(description="Remap updated_timestamps in
SIEIS data files")
    parser.add_argument("--hist-only", action="store_true", help="Remap
historical file only")
    parser.add_argument("--incr-only", action="store_true", help="Remap
incremental file only")
    parser.add_argument("--dry-run", action="store_true", help="Preview
remapping without writing")
    parser.add_argument(
        "--hist-file", default="data/processed/historical_data.txt",
        help="Path to historical file"
    )
    parser.add_argument(
        "--incr-file", default="data/processed/incremental_data.txt",
        help="Path to incremental file"
    )
    args = parser.parse_args()

    do_hist = not args.incr_only
    do_incr = not args.hist_only

    hist_path = Path(args.hist_file)
    incr_path = Path(args.incr_file)

    print("=" * 60)
    print("SIEIS Timestamp Remap")
    print(f"Today      : {date.today()}")
    print(f"Yesterday : {date.today() - timedelta(days=1)}")
    print(f"Dry run   : {args.dry_run}")
    print("=" * 60)

    results = {}

    if do_hist:
        if not hist_path.exists():
            print(f"\n✗ Historical file not found: {hist_path}")
        else:
            results["historical"] = remap_historical(hist_path,
dry_run=args.dry_run)

    if do_incr:

```

```

        if not incr_path.exists():
            print(f"\n✗ Incremental file not found: {incr_path}")
        else:
            results["incremental"] = remap_incremental(incr_path,
dry_run=args.dry_run)

# — summary


---


print("\n" + "=" * 60)
print("SUMMARY")
print("=" * 60)

if "historical" in results and results["historical"]:
    r = results["historical"]
    print(f"  Historical : {r.get('max_before')} → {r.get('max_after')} "
          f"({r.get('offset_days')} days)  {r.get('changed', 'n/a')}, {"
lines")

if "incremental" in results and results["incremental"]:
    r = results["incremental"]
    print(f"  Incremental: {r.get('unique_dates')} unique dates "
          f"→ {r.get('first_day')} ... {r.get('last_day')}")

if not args.dry_run:
    print()
    print("Next steps:")
    print("  1. Load historical to MinIO only:")
    print("      python scripts/load_historical_data.py --minio-only")
    print("  2. Start Docker pipeline for incremental stream:")
    print("      docker compose up -d")
print("=" * 60)

if __name__ == "__main__":
    main()

```

[scripts/retrain_model.py](#)

"""Retrain the SIEIS anomaly detection model with fresh data.

Usage:

```
python scripts/retrain_model.py [--days 30] [--source local|minio]
```

This script is for incremental retraining – run it after new data arrives.
Like refreshing a spam filter after the spam landscape changes.

"""

```

import sys
import os
sys.path.insert(0,

```

```

os.path.dirname(os.path.dirname(os.path.abspath(__file__))))

from dotenv import load_dotenv
load_dotenv()

import argparse
import json
import logging

logging.basicConfig(level=logging.INFO, format="% (asctime)s [% (levelname)s]
%(message)s")
logger = logging.getLogger(__name__)

def main():
    parser = argparse.ArgumentParser(description="Retrain SIEIS anomaly
detection model")
    parser.add_argument("--days", type=int, default=30, help="Days of recent
data to use")
    parser.add_argument(
        "--source",
        choices=["local", "minio", "influxdb"],
        default="local",
        help="Data source for retraining",
    )
    parser.add_argument("--contamination", type=float, default=0.05)
    args = parser.parse_args()

    from src.app import config
    from src.app.ml.detector import _load_registry, train_anomaly_detector,
save_model
    from src.app.ml.preprocessing.data_prep import (
        load_from_local_file,
        load_parquet_from_minio,
        prepare_features,
    )

    print("=" * 60)
    print("SIEIS - Incremental Model Retraining")
    print(f"Source: {args.source} | Days: {args.days}")
    print("=" * 60)

    # Show current registry state
    registry = _load_registry()
    print(f"\nCurrent registry: {len(registry.get('models', []))} model(s)")
    if registry.get("latest"):
        print(f"Current latest: {registry['latest']}")

    # Load data
    print(f"\n[1/4] Loading data...")

```



```

if args.source == "minio":
    df = load_parquet_from_minio(max_files=50)
    if df.empty:
        print("⚠️ No MinIO data, falling back to local")
        args.source = "local"

if args.source == "influxdb":
    print("Loading recent data from InfluxDB...")
    try:
        from influxdb_client import InfluxDBClient
        client = InfluxDBClient(
            url=config.INFLUX_URL, token=config.INFLUX_TOKEN,
            org=config.INFLUX_ORG
        )
        query_api = client.query_api()
        flux = f"""
from(bucket: "{config.INFLUX_BUCKET}")
  |> range(start: -{args.days}d)
  |> filter(fn: (r) => r["_measurement"] == "sensor_reading")
  |> pivot(rowKey: ["_time", "mote_id"], columnKey: ["_field"], valueColumn:
"_value")
  |> limit(n: 200000)
"""

        import pandas as pd
        tables = query_api.query_data_frame(flux, org=config.INFLUX_ORG)
        client.close()
        if isinstance(tables, list):
            df = pd.concat(tables, ignore_index=True)
        else:
            df = tables
        # Rename _time to timestamp
        if "_time" in df.columns:
            df = df.rename(columns={"_time": "timestamp"})
        print(f"    Loaded {len(df):,} rows from InfluxDB")
    except Exception as e:
        print(f"⚠️ InfluxDB load failed ({e}), falling back to local")
        args.source = "local"

if args.source == "local":
    hist_path = config.DATA_DIR / "processed" / "historical_data.txt"
    if not hist_path.exists():
        print(f"❌ historical_data.txt not found at {hist_path}")
        sys.exit(1)
    df = load_from_local_file(str(hist_path), max_rows=200_000)

print(f"    Loaded {len(df):,} rows")

# Prepare features
print("\n[2/4] Preparing features...")

```

```

X, mote_ids = prepare_features(df)
print(f"    Feature matrix: {X.shape}")

# Train
print(f"\n[3/4] Training new model version...")
pipeline, metrics = train_anomaly_detector(X,
contamination=args.contamination)
    print(f"    Anomalies: {metrics['n_anomalies_detected']:,}
({metrics['anomaly_ratio']:.1%})")

# Save with retrain tag
print("\n[4/4] Saving new model version...")
filename = save_model(pipeline, metrics, tag="retrain")

print("\n" + "=" * 60)
print("RETRAINING COMPLETE")
print("=" * 60)
print(f"    New model: {filename}")
print(f"    Reload API: curl -X POST
http://localhost:8000/api/v1/ml/model/reload")

if __name__ == "__main__":
    main()

```

scripts/run_dashboard.py

"""Launch the SIEIS Streamlit dashboard.

Usage:

```
python scripts/run_dashboard.py
"""
```

```
import os
import sys
import subprocess
```

```
# Ensure we're in the project root
project_root = os.path.dirname(os.path.dirname(os.path.abspath(__file__)))
os.chdir(project_root)
```

```
dashboard_path = os.path.join(project_root, "src", "app", "dashboard",
"app.py")
```

```
print("🚀 Starting SIEIS Dashboard...")
print(f"    Dashboard path: {dashboard_path}")
print(f"    URL: http://localhost:8501")
print("    Press Ctrl+C to stop\n")
```

```

subprocess.run([
    sys.executable, "-m", "streamlit", "run",
    dashboard_path,
    "--server.port", "8501",
    "--server.address", "0.0.0.0",
    "--browser.gatherUsageStats", "false",
], check=True)

```

scripts/simple_verify.py

```

"""Simple verification: Check if 80% bulk load succeeded."""
import sys
from pathlib import Path
sys.path.insert(0, str(Path(__file__).parent.parent))

from dotenv import load_dotenv
load_dotenv()

import pandas as pd
from src.app.config import INFLUX_URL, INFLUX_TOKEN, INFLUX_ORG,
INFLUX_BUCKET, DATA_DIR
DATA_PATH = DATA_DIR / "processed" / "historical_data.txt"
from influxdb_client import InfluxDBClient

print("🔍 SIMPLE 80% DATA LOAD VERIFICATION")
print("="*50)

# Expected 80% counts
df = pd.read_csv(DATA_PATH, sep=r'\s+', header=None,
                 names=['date', 'time', 'epoch', 'moteid', 'temperature',
                        'humidity', 'light', 'voltage'],
                 na_values=['?'])
df['timestamp'] = pd.to_datetime(df['date'] + ' ' + df['time'], format='%Y-%m-%d %H:%M:%S.%f', errors='coerce')
df = df.dropna(subset=['timestamp', 'moteid', 'temperature', 'humidity'])
df = df.sort_values('timestamp')

min_ts = df['timestamp'].min()
max_ts = df['timestamp'].max()
cutoff_ts = min_ts + (max_ts - min_ts) * 0.80
df_80 = df[df['timestamp'] <= cutoff_ts]

expected_records = len(df_80)
expected_points = expected_records * 4 # 4 fields per record

print(f"Expected from CSV (80%):")
print(f"  Records: {expected_records:,}")
print(f"  Points: {expected_points:,}")
print(f"  Date range: {min_ts.date()} to {cutoff_ts.date()}")

```




```

# Check InfluxDB
client = InfluxDBClient(url=INFLUX_URL, token=INFLUX_TOKEN, org=INFLUX_ORG,
timeout=60000)
query_api = client.query_api()

# Simple count query
result = query_api.query(f'from(bucket: "{INFLUX_BUCKET}") |> range(start:
2004-01-01, stop: 2005-01-01) |> filter(fn: (r) => r._measurement ==
"sensor_reading") |> count()')
total_points = sum(record.get_value() for table in result for record in
table.records)

print(f"\nActual in InfluxDB:")
print(f"  Points: {total_points:,}")

# Check success
success_rate = (total_points / expected_points) * 100 if expected_points > 0
else 0
print(f"\nVerification:")
print(f"  Success rate: {success_rate:.1f}%")

if success_rate >= 95:
    print("  Status:  PASSED - 80% bulk load successful!")
elif success_rate >= 80:
    print("  Status:  PARTIAL - Most data loaded")
else:
    print("  Status:  FAILED - Data load incomplete")

client.close()

```

scripts/split_dataset.py

"""Split the full sensor dataset into 80% historical and 20% incremental.

The source file is data/processed/realtime_data.txt (full dataset with remapped 2026 timestamps). Lines are split in chronological order so the 80% block represents the past and the 20% block represents "new" data.

Usage:

```

python scripts/split_dataset.py
python scripts/split_dataset.py --split 0.8
python scripts/split_dataset.py --source data/processed/realtime_data.txt
python scripts/split_dataset.py --dry-run

```

Analogy: Think of the dataset like a stack of dated newspapers. We keep the oldest 80% as the archive (historical), and the newest 20% as the fresh pile that arrives incrementally.

"""

```

import sys
import os
sys.path.insert(0,
os.path.dirname(os.path.dirname(os.path.abspath(__file__))))

import argparse
import json
import shutil
from datetime import datetime
from pathlib import Path

from dotenv import load_dotenv
load_dotenv()

def count_lines(filepath: Path) -> int:
    """Count lines in a file efficiently without loading it all into
memory."""
    count = 0
    with open(filepath, "r", encoding="utf-8", errors="replace") as f:
        for _ in f:
            count += 1
    return count

def split_dataset(
    source: Path,
    historical_out: Path,
    incremental_out: Path,
    split_ratio: float = 0.8,
    dry_run: bool = False,
):
    """Split source file into historical (split_ratio) and incremental (1-
split_ratio).

    Args:
        source: Input file (full dataset with remapped timestamps)
        historical_out: Output file for first split_ratio of lines
        incremental_out: Output file for remaining lines
        split_ratio: Fraction for historical split (default 0.8 = 80%)
        dry_run: If True, only count and report without writing
    """
    print("=" * 65)
    print("SIEIS Dataset Split")
    print(f"Source : {source}")
    print(f"Split : {split_ratio:.0%} historical / {1-split_ratio:.0%}
incremental")
    print(f"Dry run: {dry_run}")
    print("=" * 65)

```

```

# — Step 1: count total lines —————
print("\n[1/3] Counting lines in source file...")
total = count_lines(source)
historical_count = int(total * split_ratio)
incremental_count = total - historical_count

print(f"  Total lines      : {total:>12,}")
print(f"  Historical (80%) : {historical_count:>12,}")
print(f"  Incremental (20%): {incremental_count:>12,}")

if dry_run:
    print("\n✅ Dry run complete – no files written.")
    return historical_count, incremental_count

# — Step 2: backup existing files if present —————
print("\n[2/3] Backing up existing split files...")
for path in (historical_out, incremental_out):
    if path.exists():
        backup = path.with_suffix(".txt.bak")
        shutil.copy2(path, backup)
        print(f"  Backed up: {path.name} → {backup.name}")

# — Step 3: stream-write the split —————
print("\n[3/3] Writing split files...")

written_hist = 0
written_incr = 0
start = datetime.now()

with open(source, "r", encoding="utf-8", errors="replace") as src, \
    open(historical_out, "w", encoding="utf-8") as hist_f, \
    open(incremental_out, "w", encoding="utf-8") as incr_f:

    for line_num, line in enumerate(src, 1):
        if line_num <= historical_count:
            hist_f.write(line)
            written_hist += 1
        else:
            incr_f.write(line)
            written_incr += 1

    # Print progress every 10%
    if line_num % max(1, total // 10) == 0:
        pct = line_num / total * 100
        elapsed = (datetime.now() - start).total_seconds()
        rate = line_num / elapsed if elapsed > 0 else 0
        print(f"  {pct:5.0f}% – {line_num:,}/{total:,} lines
({rate:,.0f} lines/sec)")

```

```

elapsed = (datetime.now() - start).total_seconds()

# — Reset checkpoint so load_historical_data starts fresh —————
checkpoint_path = Path("data/.checkpoint_historical.json")
checkpoint_path.parent.mkdir(parents=True, exist_ok=True)
with open(checkpoint_path, "w") as f:
    json.dump({
        "last_line": 0,
        "timestamp": datetime.now().isoformat(),
        "stats": {
            "total_read": 0,
            "influx_written": 0,
            "minio_written": 0,
            "skipped": 0,
            "errors": 0,
        },
        "note": f"Reset by split_dataset.py - {split_ratio:.0%}/{1-split_ratio:.0%} split",
    }, f, indent=2)
print(f"\n 🟢 Checkpoint reset: {checkpoint_path}")

# — Summary —————
hist_size_mb = historical_out.stat().st_size / (1024 * 1024)
incr_size_mb = incremental_out.stat().st_size / (1024 * 1024)

print("\n" + "=" * 65)
print("SPLIT COMPLETE")
print("=" * 65)
print(f"   Time elapsed : {elapsed:.1f}s")
print(f"   historical_data.txt : {written_hist:>10,} lines ({hist_size_mb:.1f} MB)")
print(f"   incremental_data.txt: {written_incr:>10,} lines ({incr_size_mb:.1f} MB)")
print()
print("Next steps:")
print("  1. Load historical data into InfluxDB + MinIO:")
print("      python scripts/load_historical_data.py")
print("  2. Start Docker pipeline (streams incremental data):")
print("      docker compose up -d")
print("  3. Verify data is flowing:")
print("      python scripts/verify_influxDb.py")
print("      python scripts/verify_minio_storage.py")
print("  4. Train the ML model:")
print("      python scripts/train_model.py")
print("=" * 65)

return written_hist, written_incr

```

```

def main():
    parser = argparse.ArgumentParser(description="Split SIEIS dataset into
historical/incremental")
    parser.add_argument(
        "--source",
        default="data/processed/realtime_data.txt",
        help="Source file (default: data/processed/realtime_data.txt)",
    )
    parser.add_argument(
        "--historical-out",
        default="data/processed/historical_data.txt",
        help="Output path for historical split (default:
data/processed/historical_data.txt)",
    )
    parser.add_argument(
        "--incremental-out",
        default="data/processed/incremental_data.txt",
        help="Output path for incremental split (default:
data/processed/incremental_data.txt)",
    )
    parser.add_argument(
        "--split",
        type=float,
        default=0.8,
        metavar="RATIO",
        help="Fraction of data for historical split (default: 0.8 = 80%)",
    )
    parser.add_argument(
        "--dry-run",
        action="store_true",
        help="Count lines and report only – do not write files",
    )
    args = parser.parse_args()

    source = Path(args.source)
    if not source.exists():
        print(f"✗ Source file not found: {source}")
        sys.exit(1)

    if not (0.0 < args.split < 1.0):
        print(f"✗ --split must be between 0 and 1 (got {args.split})")
        sys.exit(1)

    split_dataset(
        source=source,
        historical_out=Path(args.historical_out),
        incremental_out=Path(args.incremental_out),
        split_ratio=args.split,

```



```

        dry_run=args.dry_run,
    )

if __name__ == "__main__":
    main()

```

scripts/train_model.py

"""Train the SIEIS anomaly detection model.

Usage:

```
python scripts/train_model.py [--source local|minio] [--max-rows 500000]
```

Steps:

1. Load data from MinIO (Parquet) or local file
2. Prepare features (clean, engineer time features)
3. Train Isolation Forest
4. Save model artifact + update registry
5. Auto-reload API + test prediction

"""

```

import sys
import os
sys.path.insert(0,
os.path.dirname(os.path.dirname(os.path.abspath(__file__))))

from dotenv import load_dotenv
load_dotenv()

import argparse
import json
import logging
import requests
from pathlib import Path

logging.basicConfig(level=logging.INFO, format="%(asctime)s [%(levelname)s]
%(message)s")
logger = logging.getLogger(__name__)

# API URL - try localhost first, fallback to container name
API_URLS = [
    "http://localhost:8000",
    "http://sieis-api:8000",
]

def get_api_url():
    """Find which API URL is reachable."""

```

```

for url in API_URLS:
    try:
        r = requests.get(f"{url}/health", timeout=3)
        if r.status_code == 200:
            return url
    except Exception:
        continue
return None

def fix_registry_paths(registry_path: Path):
    """Fix Windows absolute paths in registry to relative paths for container
    compatibility."""
    if not registry_path.exists():
        return

    with open(registry_path, "r") as f:
        registry = json.load(f)

    changed = False
    for key, info in registry.items():
        if isinstance(info, dict) and "path" in info:
            p = info["path"]
            filename = Path(p).name
            relative = f"src/app/ml/models/{filename}"
            if info["path"] != relative:
                info["path"] = relative
                changed = True
                logger.info(f"Fixed path for {key}: {p} → {relative}")

    if changed:
        with open(registry_path, "w") as f:
            json.dump(registry, f, indent=2)
        print(" ✅ Registry paths normalized for container compatibility")

def reload_api_model(api_url: str):
    """Tell the API to reload the latest model from disk."""
    try:
        r = requests.post(f"{api_url}/api/v1/ml/model/reload", timeout=10)
        if r.status_code == 200:
            print(f" ✅ API model reloaded: {r.json()}")
        else:
            print(f" ⚠️ API reload returned {r.status_code}: {r.text}")
    except Exception as e:
        print(f" ⚠️ Could not reload API model: {e}")

def test_prediction(api_url: str):

```

```

"""Send a test prediction to verify the model is working end-to-end."""
payload = {
    "temperature": 22.5,
    "humidity": 55.0,
    "light": 300,
    "voltage": 2.9
}
try:
    r = requests.post(
        f"{api_url}/api/v1/ml/predict/anomaly",
        json=payload,
        timeout=10
    )
    if r.status_code == 200:
        result = r.json()
        is_anomaly = result.get("is_anomaly", "unknown")
        score = result.get("anomaly_score", "unknown")
        print(f"    ✅ Test prediction OK – is_anomaly: {is_anomaly},
score: {score}")
    else:
        print(f"    ⚠️ Prediction returned {r.status_code}: {r.text}")
except Exception as e:
    print(f"    ⚠️ Could not test prediction: {e}")

def main():
    parser = argparse.ArgumentParser(description="Train SIEIS anomaly
detection model")
    parser.add_argument(
        "--source",
        choices=["local", "minio"],
        default="minio",
        help="Data source: 'local' (historical_data.txt) or 'minio' (Parquet
files)",
    )
    parser.add_argument("--max-rows", type=int, default=0, help="Max rows to
load (0 = all rows)")
    parser.add_argument("--contamination", type=float, default=0.05,
help="Expected anomaly fraction")
    parser.add_argument("--tag", type=str, default=None, help="Optional tag
for model filename")
    parser.add_argument("--no-reload", action="store_true", help="Skip API
reload after training")
    args = parser.parse_args()

    print("=" * 60)
    print("SIEIS – Anomaly Detection Model Training")
    max_rows_label = 'ALL' if args.max_rows == 0 else f'{args.max_rows:,}'
    print(f"Source: {args.source} | Max rows: {max_rows_label}")

```

```

print("=" * 60)

from src.app import config
from src.app.ml.preprocessing.data_prep import (
    load_from_local_file,
    load_parquet_from_minio,
    prepare_features,
)
from src.app.ml.detector import train_anomaly_detector, save_model

# — Step 1: Load data —————
print(f"\n[1/5] Loading data from {args.source}...")
if args.source == "minio":
    df = load_parquet_from_minio(max_files=0) # 0 = load all files
    if df.empty:
        print("⚠ No data in MinIO. Falling back to local file.")
        args.source = "local"
    else:
        if args.max_rows > 0 and len(df) > args.max_rows:
            df = df.sample(n=args.max_rows, random_state=42)
            print(f"    Sampled down to {args.max_rows:,} rows")

if args.source == "local":
    hist_path = config.DATA_DIR / "processed" / "historical_data.txt"
    if not hist_path.exists():
        print(f"❌ historical_data.txt not found at {hist_path}")
        sys.exit(1)
    df = load_from_local_file(str(hist_path), max_rows=args.max_rows)

print(f"    ✅ Loaded {len(df):,} rows")

# — Step 2: Prepare features —————
print(f"\n[2/5] Preparing features...")
X, mote_ids = prepare_features(df)
print(f"    ✅ Feature matrix: {X.shape} | Features: {list(X.columns)}")
print(f"    ✅ Unique motes: {len(set(mote_ids))}")

if len(X) < 100:
    print(f"❌ Too few samples for training (need at least 100)")
    sys.exit(1)

# — Step 3: Train —————
print(f"\n[3/5] Training Isolation Forest
(contamination={args.contamination})...")
pipeline, metrics = train_anomaly_detector(X,
contamination=args.contamination)
print(f"    ✅ Training complete!")
print(f"    ✅ Samples: {metrics['n_samples']:,}")

```

```

    print(f"    ✅ Anomalies detected: {metrics['n_anomalies_detected']:,}
({metrics['anomaly_ratio']:.1%})")

# — Step 4: Save + fix registry paths —————
print("\n[4/5] Saving model...")
filename = save_model(pipeline, metrics, tag=args.tag)
print(f"    ✅ Model saved: src/app/ml/models/{filename}")
fix_registry_paths(Path("src/app/ml/models/model_registry.json"))

# — Step 5: Reload API + test prediction —————
print("\n[5/5] Reloading API & testing prediction...")
if args.no_reload:
    print("    🚫 Skipped (--no-reload flag set)")
else:
    api_url = get_api_url()
    if api_url:
        print(f"    API found at: {api_url}")
        reload_api_model(api_url)
        test_prediction(api_url)
    else:
        print("    ⚠️ API not reachable at localhost:8000 or sieis-
api:8000")
        print("    Run manually: curl -X POST
http://localhost:8000/api/v1/ml/model/reload")

# — Summary —————
print("\n" + "=" * 60)
print("TRAINING COMPLETE")
print("=" * 60)
print(f"  Model file      : src/app/ml/models/{filename}")
print(f"  Source          : {args.source}")
print(f"  Samples         : {metrics['n_samples']:,}")
print(f"  Anomaly ratio   : {metrics['anomaly_ratio']:.1%}")
print(f"  Unique notes    : {len(set(mote_ids))}")
print("=" * 60)

if __name__ == "__main__":
    main()

```

scripts/validate_data_quality.py

"""Validate data quality in InfluxDB and MinIO.

Usage:

```
python scripts/validate_data_quality.py
```

Checks:

- InfluxDB connectivity and recent data

```

- MinIO bucket and Parquet file count
- Data completeness (missing fields, null rates)
- Timestamp freshness
"""

import sys
import os
sys.path.insert(0,
os.path.dirname(os.path.dirname(os.path.abspath(__file__))))

from dotenv import load_dotenv
load_dotenv()

import logging
from datetime import datetime, timezone

logging.basicConfig(level=logging.INFO, format="% (asctime)s [% (levelname)s]
% (message)s")
logger = logging.getLogger(__name__)

from src.app import config

def check_influxdb() -> bool:
    """Validate InfluxDB connectivity and recent data presence."""
    print("\n🇩🇪 Checking InfluxDB...")
    try:
        from influxdb_client import InfluxDBClient
        client = InfluxDBClient(url=config.INFLUX_URL,
token=config.INFLUX_TOKEN, org=config.INFLUX_ORG)


        # Ping
        ok = client.ping()
        print(f"✅ Ping: {'OK' if ok else '❌ FAILED'}")
        if not ok:
            return False

        # Count recent records
        query_api = client.query_api()
        flux = f"""
from(bucket: "{config.INFLUX_BUCKET}")
  |> range(start: -24h)
  |> filter(fn: (r) => r["_measurement"] == "sensor_reading")
  |> filter(fn: (r) => r["_field"] == "temperature")
  |> count()
"""
        tables = query_api.query(flux, org=config.INFLUX_ORG)
        total = sum(r.get_value() or 0 for t in tables for r in t.records)
        print(f"✅ Records in last 24h: {total:,}")
    
```


```


        # Check mote count
        flux2 = f"""
from(bucket: "{config.INFLUX_BUCKET}")
  |> range(start: -24h)
  |> filter(fn: (r) => r["_measurement"] == "sensor_reading")
  |> keep(columns: ["mote_id"])
  |> distinct(column: "mote_id")
  |> count()
"""


        tables2 = query_api.query(flux2, org=config.INFLUX_ORG)
        mote_count = sum(r.get_value() or 0 for t in tables2 for r in
t.records)


        print(f"  Active motes (24h): {mote_count}")

        client.close()
        return total > 0

    except Exception as e:
        print(f"  InfluxDB check failed: {e}")
        return False

def check_minio() -> bool:
    """Validate MinIO connectivity and Parquet files."""
    print(f"\n Checking MinIO...")
    try:
        from minio import Minio
        client = Minio(
            config.MINIO_ENDPOINT,
            access_key=config.MINIO_ACCESS_KEY,
            secret_key=config.MINIO_SECRET_KEY,
            secure=config.MINIO_SECURE,
        )

        # List buckets
        buckets = client.list_buckets()
        bucket_names = [b.name for b in buckets]
        print(f"  Buckets found: {bucket_names}")

        if config.MINIO_BUCKET not in bucket_names:
            print(f"  Expected bucket '{config.MINIO_BUCKET}' not found!")
            return False

        # Count Parquet files
        objects = list(client.list_objects(config.MINIO_BUCKET,
recursive=True))
        parquet_files = [o for o in objects if
o.object_name.endswith(".parquet")]

```

```

total_size_mb = sum(o.size or 0 for o in parquet_files) / (1024 *
1024)

print(f" ✅ Parquet files: {len(parquet_files)} ({total_size_mb:.1f}
MB) ")

if parquet_files:
    sample_names = [o.object_name for o in parquet_files[:3]]
    print(f" 📁 Sample paths: {sample_names}")

    return len(parquet_files) > 0

except Exception as e:
    print(f" ❌ MinIO check failed: {e}")
    return False

def check_data_quality() -> bool:
    """Check data completeness from local historical file if InfluxDB is
empty."""
    print("\n📁 Checking local data files...")
    try:
        import pandas as pd

        hist_path = config.DATA_DIR / "processed" / "historical_data.txt"
        if not hist_path.exists():
            print(f" ⚠️ historical_data.txt not found at {hist_path}")
            return False

        cols = ["date", "time", "epoch", "mote_id", "temperature", "humidity",
"light", "voltage", "updated_timestamp"]
        df = pd.read_csv(hist_path, sep=r"\s+", names=cols, na_values=["N/A",
"nan", ""], nrows=10000, on_bad_lines="skip")

        total = len(df)
        print(f" ✅ Rows sampled: {total:,}")

        for col in ["temperature", "humidity", "light", "voltage"]:
            null_pct = df[col].isna().mean() * 100
            status = "✅" if null_pct < 20 else "⚠️ "
            print(f" {status} {col}: {null_pct:.1f}% null")

        mote_count = df["mote_id"].nunique()
        print(f" ✅ Unique motes in sample: {mote_count}")

        return True

    except Exception as e:

```



```

        print(f" ❌ Data quality check failed: {e}")
        return False

def main():
    print("=" * 60)
    print("SIEIS Data Quality Validation")
    print(f"Time: {datetime.now(timezone.utc).isoformat()}")
    print("=" * 60)

    influx_ok = check_influxdb()
    minio_ok = check_minio()
    local_ok = check_data_quality()

    print("\n" + "=" * 60)
    print("SUMMARY")
    print("=" * 60)
    print(f" InfluxDB: {'✅ OK' if influx_ok else '❌ FAILED or empty'}")
    print(f" MinIO:      {'✅ OK' if minio_ok else '❌ FAILED or empty'}")
    print(f" Local:       {'✅ OK' if local_ok else '❌ FAILED'}")

    all_ok = local_ok # local is minimum requirement
    if not all_ok:
        print("\n❌ Validation failed – run load_historical_data.py first")
        sys.exit(1)
    else:
        print("\n✅ Validation passed – ready for ML training")
        print(" Next step: python scripts/train_model.py")

if __name__ == "__main__":
    main()

```

scripts/verify_influxDb.py

```
"""
Verify InfluxDB has data in the sensor_data bucket.
"""
from influxdb_client import InfluxDBClient
import logging
import sys

# Force UTF-8 output to avoid encoding errors
if hasattr(sys.stdout, 'reconfigure'):
    sys.stdout.reconfigure(encoding='utf-8')

logging.basicConfig(level=logging.INFO, format='%(levelname)s: %(message)s')
logger = logging.getLogger("InfluxDB-Verify")

INFLUXDB_URL = "http://localhost:8086"
INFLUXDB_TOKEN = "my-super-secret-token"
INFLUXDB_ORG = "sieis"
INFLUXDB_BUCKET = "sensor_data"

def check_date_range(client):
    """Check the earliest and latest timestamps in the bucket."""
    query_api = client.query_api()

    # Query for earliest record
    earliest_query = f'''
    from(bucket:"{INFLUXDB_BUCKET}")
      |> range(start: 0)
      |> sort(columns: ["_time"])
      |> limit(n:1)
    '''

    # Query for latest record
    latest_query = f'''
    from(bucket:"{INFLUXDB_BUCKET}")
      |> range(start: 0)
      |> sort(columns: ["_time"], desc: true)
      |> limit(n:1)
    '''

    try:
        logger.info("Checking date range of stored data...")

        earliest_tables = query_api.query(earliest_query)
        latest_tables = query_api.query(latest_query)

        earliest_time = None
        latest_time = None
```

```

for table in earliest_tables:
    if table.records:
        earliest_time = table.records[0].get_time()
        break

for table in latest_tables:
    if table.records:
        latest_time = table.records[0].get_time()
        break

if earliest_time and latest_time:
    logger.info(f"📅 Date Range:")
    logger.info(f"    Earliest: {earliest_time}")
    logger.info(f"    Latest:    {latest_time}")
    logger.info(f"    Span:      {latest_time - earliest_time}")
    return earliest_time, latest_time
else:
    logger.warning(f"⚠️ Could not determine date range (no data
found)")
    return None, None

except Exception as e:
    logger.error(f"❌ Error checking date range: {e}")
    return None, None

def main():
    logger.info(f"Connecting to InfluxDB at {INFLUXDB_URL}...")

    try:
        client = InfluxDBClient(url=INFLUXDB_URL, token=INFLUXDB_TOKEN,
org=INFLUXDB_ORG)
        query_api = client.query_api()

        # First check date range
        earliest, latest = check_date_range(client)

        if not earliest:
            logger.warning(f"⚠️ No data found in InfluxDB at all!")
            logger.info("Possible reasons:")
            logger.info("  1) Consumer hasn't started yet")
            logger.info("  2) Simulator hasn't produced data")
            logger.info("  3) Check: docker logs sieis-consumer")
            client.close()
            return

        # Query for ALL data first (no measurement filter) to see what's there
        all_data_query = f'''
from(bucket:"{INFLUXDB_BUCKET}")
  |> range(start: 0)

```

```

    |> limit(n:1000)
'''

logger.info("Querying all data in bucket...")
all_tables = query_api.query(all_data_query)

# Collect measurements, fields, and other info
total_count = 0
mote_ids = set()
fields = set()
measurements = set()

for table in all_tables:
    for record in table.records:
        total_count += 1
        mote_id = record.values.get("mote_id")
        field = record.get_field()
        measurement = record.get_measurement()

        if mote_id:
            mote_ids.add(mote_id)
        if field:
            fields.add(field)
        if measurement:
            measurements.add(measurement)

logger.info(f"✅ Total records (sample of 1000): {total_count:,}")
logger.info(f"✅ Measurements found: {measurements}")
logger.info(f"✅ Fields found: {fields}")

# Get unique motes from actual data
logger.info(f"✅ Unique motes: {len(mote_ids)}")
if mote_ids:
    sample_motes = sorted(list(mote_ids))[:10]
    logger.info(f"    Sample: {sample_motes}")

# Query for sample recent data
query = f'''
from(bucket:"{INFLUXDB_BUCKET}")
  |> range(start: -24h)
  |> filter(fn: (r) => r._measurement == "sensor_reading")
  |> limit(n:5)
'''

logger.info("")
logger.info("Querying for sample records (last 24h)...")
tables = query_api.query(query)

record_count = sum(len(table.records) for table in tables)

```

```

        if record_count > 0:
            logger.info(f"✅ Found {record_count} recent records")
            logger.info("Sample records:")
            for table in tables[:1]:
                for record in table.records[:3]:
                    logger.info(f"    {record.get_time()}")
mote={record.values.get('mote_id')}
{record.get_field()}={record.get_value()}")
        else:
            logger.warning("⚠️ No records in last 24h (but older data
exists)")

        logger.info("")
        logger.info(f"✅ InfluxDB verification complete!")
        client.close()

    except Exception as e:
        logger.error(f"❌ Error querying InfluxDB: {e}")
        import traceback
        traceback.print_exc()

if __name__ == "__main__":
    main()

```

scripts/verify_minio_storage.py

```

"""
Verify MinIO storage has expected files.
List objects in 'sieis-archive' and check if they match today's date.
"""

from minio import Minio
import logging
from datetime import datetime

# Configure logging
logging.basicConfig(level=logging.INFO, format='%(asctime)s - %(message)s')
logger = logging.getLogger("MinIO-Verify")

# Config matches docker-compose
MINIO_ENDPOINT = "localhost:9000"
MINIO_ACCESS_KEY = "minioadmin"
MINIO_SECRET_KEY = "minioadmin123"
BUCKET_NAME = "sieis-archive"

def main():
    logger.info(f"Connecting to MinIO at {MINIO_ENDPOINT}...")

    try:

```

```

client = Minio(
    MINIO_ENDPOINT,
    access_key=MINIO_ACCESS_KEY,
    secret_key=MINIO_SECRET_KEY,
    secure=False
)

if not client.bucket_exists(BUCKET_NAME):
    logger.error(f"✗ Bucket '{BUCKET_NAME}' does not exist!")
    return

logger.info(f"Checking bucket '{BUCKET_NAME}'...")
objects = list(client.list_objects(BUCKET_NAME, recursive=True))

count = len(objects)
logger.info(f"✅ Found {count} objects in bucket.")

if count == 0:
    logger.warning("    Bucket is empty. Consumer might not have
written data yet.")
    return

# Check for files matching today's date (formatted as YYYYMMDD in our
parquet script)
today_str = datetime.now().strftime("%Y%m%d")
today_files = [o for o in objects if today_str in o.object_name]

if today_files:
    logger.info(f"✅ Found {len(today_files)} files matching today's
date pattern ({today_str}).")
    logger.info("Sample files:")
    for o in today_files[:3]:
        logger.info(f"    - {o.object_name} ({o.size} bytes)")
else:
    logger.warning(f"⚠️ No files found matching today's date pattern
({today_str}).")
    logger.info("Sample of existing files:")
    for o in objects[:3]:
        logger.info(f"    - {o.object_name}")

except Exception as e:
    logger.error(f"✗ Error verifying MinIO: {e}")

if __name__ == "__main__":
    main()

```

src/__init__.py

```
"""Top-level package for SIEIS."""
```

src/app/__init__.py

```
"""Application package for SIEIS."""
```

src/app/api/__init__.py

src/app/api/main.py

```
"""SIEIS FastAPI application - main entry point."""
```

```
import logging
```

```
import os
```

```
from contextlib import asynccontextmanager
```

```
from fastapi import FastAPI
```

```
from fastapi.middleware.cors import CORSMiddleware
```

```
from fastapi.responses import JSONResponse
```

```
from src.app.api.routes import sensors, analytics, ml
```

```
from src.app import config
```

```
logging.basicConfig(
```

```
    level=getattr(logging, config.LOG_LEVEL, logging.INFO),
```

```
    format="%(asctime)s [%(levelname)s] %(name)s: %(message)s",
```

```
)
```

```
logger = logging.getLogger(__name__)
```

```
@asynccontextmanager
```

```
async def lifespan(app: FastAPI):
```

```
    """Startup and shutdown events."""
```

```
    logger.info(f"🚀 SIEIS API starting - env={config.ENVIRONMENT},
```

```
influxdb={config.INFLUX_URL}")
```

```
    yield
```

```
    logger.info("SIEIS API shutting down")
```

```
app = FastAPI(
```

```
    title=config.API_TITLE,
```

```
    version=config.API_VERSION,
```

```
    description="SIEIS - Smart IoT Environmental Information System API",
```

```
    lifespan=lifespan,
```

```
)
```

```

# CORS - allow all origins for local development
app.add_middleware(
    CORSMiddleware,
    allow_origins=["*"],
    allow_credentials=True,
    allow_methods=["*"],
    allow_headers=["*"],
)

# Register routers
app.include_router(sensors.router, prefix="/api/v1")
app.include_router(analytics.router, prefix="/api/v1")
app.include_router(ml.router, prefix="/api/v1")

@app.get("/", tags=["root"])
def root():
    return {"message": "SIEIS API", "version": config.API_VERSION, "docs":
"/docs"}

@app.get("/health", tags=["health"])
@app.get("/api/v1/health", tags=["health"])
def health():
    """Basic health check - verifies InfluxDB connectivity."""
    influx_status = "unknown"
    try:
        from influxdb_client import InfluxDBClient
        client = InfluxDBClient(
            url=config.INFLUX_URL,
            token=config.INFLUX_TOKEN,
            org=config.INFLUX_ORG,
        )
        ping_ok = client.ping()
        influx_status = "ok" if ping_ok else "error"
        client.close()
    except Exception as e:
        influx_status = f"error: {e}"

    # Check if model is loaded
    from src.app.api.routes.ml import _get_model
    model_loaded = _get_model() is not None

    status = "ok" if influx_status == "ok" else "degraded"
    return JSONResponse(
        status_code=200,
        content={
            "status": status,
            "influxdb": influx_status,
            "model_loaded": model_loaded,

```



```

        "version": config.API_VERSION,
    },
)

```

`src/app/api/routes/__init__.py`

`src/app/api/routes/analytics.py`

```

"""Analytics endpoints - aggregations, trends, statistics."""

import logging
from typing import List, Optional

from fastapi import APIRouter, HTTPException, Query

from src.app.api.schemas import AggregatedStats
from src.app import config

logger = logging.getLogger(__name__)
router = APIRouter(prefix="/analytics", tags=["analytics"])

def _get_influx_client():
    from influxdb_client import InfluxDBClient
    return InfluxDBClient(
        url=config.INFLUX_URL,
        token=config.INFLUX_TOKEN,
        org=config.INFLUX_ORG,
    )

@router.get("/summary", summary="Get aggregated stats for a metric")
def get_summary(
    metric: str = Query(..., description="Sensor metric: temperature, humidity, light, voltage"),
    mote_id: Optional[str] = Query(default=None, description="Filter by mote ID (optional)"),
    start: str = Query(default="-24h", description="Start time"),
):
    """Return mean, min, max for a given sensor metric over a time window."""
    valid_metrics = {"temperature", "humidity", "light", "voltage"}
    if metric not in valid_metrics:
        raise HTTPException(status_code=400, detail=f"metric must be one of {valid_metrics}")

    try:
        client = _get_influx_client()
        query_api = client.query_api()

```

```

        mote_filter = f'|> filter(fn: (r) => r["mote_id"] == "{mote_id}")' if
mote_id else ""

        flux = f"""
from(bucket: "{config.INFLUX_BUCKET}")
  |> range(start: {start})
  |> filter(fn: (r) => r["_measurement"] == "sensor_reading")
  |> filter(fn: (r) => r["_field"] == "{metric}")
  {mote_filter}
"""

        tables = query_api.query(flux, org=config.INFLUX_ORG)

        values = []
        for table in tables:
            for record in table.records:
                v = record.get_value()
                if v is not None:
                    values.append(float(v))

        client.close()

        if not values:
            return AggregatedStats(
                mote_id=mote_id, metric=metric,
                mean=None, min=None, max=None, count=0, period=start
            )

        return AggregatedStats(
            mote_id=mote_id,
            metric=metric,
            mean=round(sum(values) / len(values), 4),
            min=round(min(values), 4),
            max=round(max(values), 4),
            count=len(values),
            period=start,
        )
    except Exception as e:
        logger.exception("Failed to compute summary")
        raise HTTPException(status_code=500, detail=str(e))

@router.get("/active-motes-count", summary="Count active motes in last 15
minutes")
def active_motes_count():
    """Return count of motes that reported data in the last 15 minutes."""
    try:
        client = _get_influx_client()
        query_api = client.query_api()

        flux = f"""

```

```

from(bucket: "{config.INFLUX_BUCKET}")
  |> range(start: -15m)
  |> filter(fn: (r) => r["_measurement"] == "sensor_reading")
  |> keep(columns: ["mote_id"])
  |> distinct(column: "mote_id")
  |> count()
"""
    tables = query_api.query(flux, org=config.INFLUX_ORG)
    client.close()

    count = 0
    for table in tables:
        for record in table.records:
            count = record.get_value() or 0
    return {"active_motes_last_15m": int(count)}
except Exception as e:
    logger.exception("Failed to count active motes")
    raise HTTPException(status_code=500, detail=str(e))

@router.get("/ingestion-rate", summary="Messages ingested per minute (last 5m)")
def ingestion_rate():
    """Return approximate records-per-minute ingestion rate."""
    try:
        client = _get_influx_client()
        query_api = client.query_api()

        flux = f"""
from(bucket: "{config.INFLUX_BUCKET}")
  |> range(start: -5m)
  |> filter(fn: (r) => r["_measurement"] == "sensor_reading")
  |> filter(fn: (r) => r["_field"] == "temperature")
  |> count()
"""
        tables = query_api.query(flux, org=config.INFLUX_ORG)
        client.close()

        total = 0
        for table in tables:
            for record in table.records:
                total += record.get_value() or 0
        rate = round(total / 5.0, 2)
        return {"records_per_minute": rate, "total_last_5m": total}
    except Exception as e:
        logger.exception("Failed to compute ingestion rate")
        raise HTTPException(status_code=500, detail=str(e))

```

`src/app/api/routes/ml.py`

```
"""ML inference endpoints - anomaly detection."""

import logging
import os
from typing import List

from fastapi import APIRouter, HTTPException

from src.app.api.schemas import AnomalyInput, AnomalyResult
from src.app import config

logger = logging.getLogger(__name__)
router = APIRouter(prefix="/ml", tags=["ml"])

# Global model cache (loaded once on first use)
_model = None
_model_path = None

def _load_model():
    """Load the latest anomaly detection model from models directory."""
    global _model, _model_path
    import pickle
    import json

    registry_path = config.MODELS_DIR / "model_registry.json"
    if not registry_path.exists():
        return None, None

    with open(registry_path) as f:
        registry = json.load(f)

    latest = registry.get("latest")
    if not latest:
        return None, None

    model_file = config.MODELS_DIR / latest
    if not model_file.exists():
        return None, None

    with open(model_file, "rb") as f:
        artifact = pickle.load(f)

    # Artifact is a dict {"pipeline": ..., "metrics": ..., "version": ...}
    if isinstance(artifact, dict) and "pipeline" in artifact:
        model = artifact["pipeline"]
    else:
        model = artifact # backwards compat for plain pipeline pickles
```

```

    return model, str(model_file)

def _get_model():
    global _model, _model_path
    if _model is None:
        _model, _model_path = _load_model()
    return _model

def _score_to_severity(score: float) -> str:
    """Map anomaly score to human-readable severity.

    The anomaly_score is computed as: 0.5 - decision_function_score / 2
    So the natural split points are:
        < 0.50 → IsolationForest decision_function > 0 → model says normal
        0.50-0.65 → mild anomaly
        >= 0.65 → clear anomaly
    """
    if score >= 0.65:
        return "critical"
    elif score >= 0.5:
        return "warning"
    return "normal"

@router.post("/predict/anomaly", response_model=AnomalyResult,
summary="Predict if a reading is anomalous")
def predict_anomaly(body: AnomalyInput):
    """Run anomaly detection on a single sensor reading.

    Returns an anomaly score (0=normal, 1=anomaly) and severity label.
    Falls back to a rule-based heuristic if no model is loaded.
    """
    model = _get_model()

    if model is not None:
        import numpy as np
        import pandas as pd
        from datetime import datetime

        # Build feature vector matching training: temperature, humidity,
light, voltage, hour, day_of_week
        now = datetime.utcnow()
        X = pd.DataFrame([
            "temperature": body.temperature,
            "humidity": body.humidity,
            "light": body.light,
            "voltage": body.voltage,

```

```

        "hour": now.hour,
        "day_of_week": now.weekday(),
    })
    # Isolation Forest: predict returns -1 (anomaly) or 1 (normal)
    raw_pred = model.predict(X)[0]
    score_arr = model.decision_function(X)[0]
    # Convert decision function score to 0-1 range (higher = more
anomalous)
    anomaly_score = float(max(0.0, min(1.0, 0.5 - score_arr / 2.0)))
    is_anomaly = raw_pred == -1
else:
    # Rule-based fallback when no model is trained yet
    anomaly_score = 0.0
    flags = []
    if not (15 <= body.temperature <= 45):
        flags.append("temperature"); anomaly_score += 0.4
    if not (20 <= body.humidity <= 90):
        flags.append("humidity"); anomaly_score += 0.3
    if body.light < 0 or body.light > 2000:
        flags.append("light"); anomaly_score += 0.2
    if not (2.0 <= body.voltage <= 3.5):
        flags.append("voltage"); anomaly_score += 0.3
    anomaly_score = min(1.0, anomaly_score)
    is_anomaly = anomaly_score > 0.3

return AnomalyResult(
    mote_id=body.mote_id,
    anomaly_score=round(anomaly_score, 4),
    is_anomaly=is_anomaly,
    severity=_score_to_severity(anomaly_score),
    input=body,
)

```

```

@router.get("/model/info", summary="Get info about the currently loaded
model")

```

```

def model_info():
    """Return metadata about the loaded anomaly detection model."""
    import json
    registry_path = config.MODELS_DIR / "model_registry.json"
    if registry_path.exists():
        with open(registry_path) as f:
            registry = json.load(f)
    else:
        registry = {}

    model = _get_model()
    return {
        "model_loaded": model is not None,
        "model_path": _model_path,
    }

```

```

        "registry": registry,
        "fallback_mode": model is None,
    }

@router.post("/model/reload", summary="Reload model from disk")
def reload_model():
    """Force-reload the latest model from the models directory."""
    global _model, _model_path
    _model, _model_path = _load_model()
    return {
        "reloaded": True,
        "model_loaded": _model is not None,
        "model_path": _model_path,
    }

```

[src/app/api/routes/sensors.py](#)

```

"""Sensor data query endpoints - reads from InfluxDB."""

import logging
from typing import List, Optional

from fastapi import APIRouter, HTTPException, Query

from src.app.api.schemas import LatestReading, SensorReading
from src.app import config

logger = logging.getLogger(__name__)
router = APIRouter(prefix="/sensors", tags=["sensors"])

def _get_influx_client():
    from influxdb_client import InfluxDBClient
    return InfluxDBClient(
        url=config.INFLUX_URL,
        token=config.INFLUX_TOKEN,
        org=config.INFLUX_ORG,
    )

@router.get("/latest", response_model=List[LatestReading], summary="Get latest reading per mote")
def get_latest_readings(
    mote_id: Optional[str] = Query(default=None, description="Filter by mote ID"),
    window: str = Query(default="-1h", description="Time window, e.g. -1h, -15m, -24h"),
):

```

```

    """Return the most recent sensor reading for each mote (or a specific
mote)."""
    try:
        client = _get_influx_client()
        query_api = client.query_api()

        mote_filter = f'|> filter(fn: (r) => r["mote_id"] == "{mote_id}")' if
mote_id else ""

        flux = f"""
from(bucket: "{config.INFLUX_BUCKET}")
  |> range(start: {window})
  |> filter(fn: (r) => r["_measurement"] == "sensor_reading")
  {mote_filter}
  |> last()
  |> pivot(rowKey:["_time","mote_id"], columnKey: ["_field"], valueColumn:
"_value")
"""

        tables = query_api.query(flux, org=config.INFLUX_ORG)
        client.close()

        results = []
        for table in tables:
            for record in table.records:
                results.append(LatestReading(
                    mote_id=record.values.get("mote_id", "unknown"),
                    temperature=record.values.get("temperature"),
                    humidity=record.values.get("humidity"),
                    light=record.values.get("light"),
                    voltage=record.values.get("voltage"),
                    timestamp=str(record.get_time()) if record.get_time() else
None,

                    ))
            return results
    except Exception as e:
        logger.exception("Failed to query latest readings")
        raise HTTPException(status_code=500, detail=str(e))

@router.get("/history", response_model=List[SensorReading], summary="Get
sensor readings over time")
def get_sensor_history(
    mote_id: str = Query(..., description="Mote ID to query"),
    start: str = Query(default="-1h", description="Start time, e.g. -1h, -24h,
or ISO timestamp"),
    stop: str = Query(default="now()", description="Stop time"),
    limit: int = Query(default=500, ge=1, le=5000),
):
    """Return time-series sensor readings for a specific mote."""
    try:

```



```

        client = _get_influx_client()
        query_api = client.query_api()

        flux = f"""
from(bucket: "{config.INFLUX_BUCKET}")
  |> range(start: {start}, stop: {stop})
  |> filter(fn: (r) => r["_measurement"] == "sensor_reading")
  |> filter(fn: (r) => r["mote_id"] == "{mote_id}")
  |> pivot(rowKey: ["_time", "mote_id"], columnKey: ["_field"], valueColumn:
"_value")
  |> limit(n: {limit})
"""

        tables = query_api.query(flux, org=config.INFLUX_ORG)
        client.close()

        results = []
        for table in tables:
            for record in table.records:
                results.append(SensorReading(
                    mote_id=record.values.get("mote_id", mote_id),
                    temperature=record.values.get("temperature"),
                    humidity=record.values.get("humidity"),
                    light=record.values.get("light"),
                    voltage=record.values.get("voltage"),
                    timestamp=record.get_time(),
                ))
        return results
    except Exception as e:
        logger.exception("Failed to query sensor history")
        raise HTTPException(status_code=500, detail=str(e))

@router.get("/motes", response_model=List[str], summary="List all active mote
IDs")
def list_motes(
    window: str = Query(default="-24h", description="Time window to look for
active motes"),
):
    """Return list of mote IDs that have data in the given time window."""
    try:
        client = _get_influx_client()
        query_api = client.query_api()

        flux = f"""
from(bucket: "{config.INFLUX_BUCKET}")
  |> range(start: {window})
  |> filter(fn: (r) => r["_measurement"] == "sensor_reading")
  |> keep(columns: ["mote_id"])
  |> distinct(column: "mote_id")
"""

```

```

tables = query_api.query(flux, org=config.INFLUX_ORG)
client.close()

mote_ids = []
for table in tables:
    for record in table.records:
        val = record.values.get("mote_id") or record.get_value()
        if val and val not in mote_ids:
            mote_ids.append(str(val))
    return sorted(mote_ids)
except Exception as e:
    logger.exception("Failed to list motes")
    raise HTTPException(status_code=500, detail=str(e))

```

[src/app/api/schemas.py](#)

"""Pydantic schemas for SIEIS API request/response models."""

```

from datetime import datetime
from typing import List, Optional
from pydantic import BaseModel, Field

```

```

class SensorReading(BaseModel):
    mote_id: str
    temperature: Optional[float] = None
    humidity: Optional[float] = None
    light: Optional[float] = None
    voltage: Optional[float] = None
    timestamp: Optional[datetime] = None

```

```

class SensorQuery(BaseModel):
    mote_id: Optional[str] = None
    start: str = Field(default="-1h", description="Flux duration or ISO
timestamp, e.g. '-1h', '-24h', '2026-01-01T00:00:00Z'")
    stop: str = Field(default="now()", description="Flux stop time")
    limit: int = Field(default=500, ge=1, le=10000)

```

```

class LatestReading(BaseModel):
    mote_id: str
    temperature: Optional[float] = None
    humidity: Optional[float] = None
    light: Optional[float] = None
    voltage: Optional[float] = None
    timestamp: Optional[str] = None

```

```

class AggregatedStats(BaseModel):
    mote_id: Optional[str] = None
    metric: str
    mean: Optional[float] = None
    min: Optional[float] = None
    max: Optional[float] = None
    count: Optional[int] = None
    period: str

class AnomalyInput(BaseModel):
    temperature: float
    humidity: float
    light: float
    voltage: float
    mote_id: Optional[str] = None

class AnomalyResult(BaseModel):
    mote_id: Optional[str] = None
    anomaly_score: float
    is_anomaly: bool
    severity: str # "normal", "warning", "critical"
    input: AnomalyInput

class HealthResponse(BaseModel):
    status: str
    influxdb: str
    model_loaded: bool
    version: str = "1.0.0"

```

[src/app/api_server.py](#)

"""Entry point to run the SIEIS API server.

Usage:

```
python -m src.app.api_server
```

Or with uvicorn directly:

```
uvicorn src.app.api.main:app --host 0.0.0.0 --port 8000 --reload
```

"""

```
import uvicorn
```

```
from dotenv import load_dotenv
```

```
# Load .env before importing config
```

```
load_dotenv()
```

```

from src.app import config

if __name__ == "__main__":
    uvicorn.run(
        "src.app.api.main:app",
        host=config.API_HOST,
        port=config.API_PORT,
        reload=config.DEBUG,
        log_level=config.LOG_LEVEL.lower(),
    )

```

src/app/config.py

```

import os
from pathlib import Path

# =====
# PROJECT PATHS
# =====
PROJECT_ROOT = Path(__file__).parent.parent.parent
DATA_DIR = PROJECT_ROOT / "data"
LOGS_DIR = PROJECT_ROOT / "logs"
MODELS_DIR = PROJECT_ROOT / "src/app/ml/models"

# =====
# INFLUXDB CONFIGURATION
# =====
INFLUX_URL = os.getenv("INFLUX_URL", "http://localhost:8086")
INFLUX_TOKEN = os.getenv("INFLUX_TOKEN", "my-super-secret-token")
INFLUX_ORG = os.getenv("INFLUX_ORG", "sieis")
INFLUX_BUCKET = os.getenv("INFLUX_BUCKET", "sensor_data")

# =====
# MINIO CONFIGURATION (S3-COMPATIBLE OBJECT STORAGE)
# =====
MINIO_ENDPOINT = os.getenv("MINIO_ENDPOINT", "localhost:9000")
MINIO_ACCESS_KEY = os.getenv("MINIO_ACCESS_KEY", "minioadmin")
MINIO_SECRET_KEY = os.getenv("MINIO_SECRET_KEY", "minioadmin123")
MINIO_SECURE = os.getenv("MINIO_SECURE", "false").lower() == "true"
MINIO_BUCKET = os.getenv("MINIO_BUCKET", "sieis-archive")

# =====
# REDPANDA/KAFKA CONFIGURATION
# =====
KAFKA_BROKER = os.getenv("KAFKA_BROKER", "localhost:19092") # 29092 was
wrong; Redpanda external port is 19092
KAFKA_TOPIC = os.getenv("KAFKA_TOPIC", "sensor_readings")
KAFKA_GROUP = os.getenv("KAFKA_GROUP", "sieis-consumer-group")

```

```

# =====
# SIMULATOR CONFIGURATION
# =====
SPEED_FACTOR = float(os.getenv("SPEED_FACTOR", "100.0"))
DATA_PATH = os.getenv("DATA_PATH", str(DATA_DIR /
"processed/incremental_data.txt"))
MOTE_LOCS_PATH = os.getenv("MOTE_LOCS_PATH", str(DATA_DIR /
"raw/mote_locs.txt"))
FILTER_TODAY_ONLY = os.getenv("FILTER_TODAY_ONLY", "true").lower() == "true"
# YEAR_OFFSET: shifts legacy 2004 timestamps to current year (fallback when
updated_timestamp absent)
import datetime as _dt
YEAR_OFFSET = int(os.getenv("YEAR_OFFSET", str(_dt.datetime.now().year -
2004)))

# =====
# API CONFIGURATION
# =====
API_HOST = os.getenv("API_HOST", "0.0.0.0")
API_PORT = int(os.getenv("API_PORT", "8000"))
API_TITLE = "SIEIS API"
API_VERSION = "1.0.0"

# =====
# LOGGING CONFIGURATION
# =====
LOG_LEVEL = os.getenv("LOG_LEVEL", "INFO")
LOG_FILE = LOGS_DIR / "sieis.log"

# =====
# APPLICATION CONFIGURATION
# =====
DEBUG = os.getenv("DEBUG", "false").lower() == "true"
ENVIRONMENT = os.getenv("ENVIRONMENT", "development")

# =====
# VALIDATION
# =====
def validate_config():
    """Validate all required configuration is set."""
    required = [
        ("INFLUX_URL", INFLUX_URL),
        ("MINIO_ENDPOINT", MINIO_ENDPOINT),
        ("MINIO_ACCESS_KEY", MINIO_ACCESS_KEY),
        ("MINIO_SECRET_KEY", MINIO_SECRET_KEY),
    ]

    missing = [name for name, value in required if not value]

    if missing:

```

```
        raise ValueError(f"Missing required configuration: {'',
'.join(missing)}")
```

```
if __name__ == "__main__":
    validate_config()
    print("✅ Config validation passed!")
```

[src/app/consumer/__init__.py](#)

```
"""Consumer package for SIEIS."""
```

[src/app/consumer/influx_writer.py](#)

```
"""InfluxDB 2.x writer for sensor readings."""
```

```
import logging
from datetime import datetime
from typing import Dict, Iterable, List, Optional

from influxdb_client import InfluxDBClient, Point, WritePrecision
from influxdb_client.client.write_api import SYNCHRONOUS

logger = logging.getLogger(__name__)
```

```
class InfluxWriter:
    """Write batches of sensor readings to InfluxDB 2.x."""

    def __init__(
        self,
        url: str,
        token: str,
        org: str,
        bucket: str,
        measurement: str = "sensor_reading",
    ) -> None:
        """Initialize InfluxDB 2.x writer.

        Args:
            url: InfluxDB URL (e.g., 'http://localhost:8086')
            token: Authentication token
            org: Organization name
            bucket: Bucket name
            measurement: Measurement name for time series data
        """
        self.org = org
        self.bucket = bucket
        self.measurement = measurement
        self.client = InfluxDBClient(url=url, token=token, org=org)
```

```

        self.write_api = self.client.write_api(write_options=SYNCHRONOUS)
        logger.info(f"InfluxDB 2.x writer initialized: url={url}, org={org},
bucket={bucket}")

def _to_point(self, message: Dict) -> Optional[Point]:
    """Convert message dict to InfluxDB Point.

    Uses 'updated_timestamp' for time field (maps 2004 data to 2025+).
    """
    mote_id = message.get("mote_id")
    if mote_id is None:
        logger.warning("Message missing mote_id, skipping")
        return None

    point = Point(self.measurement).tag("mote_id", str(mote_id))

    # Add sensor fields
    field_count = 0
    for field in ("temperature", "humidity", "light", "voltage"):
        value = message.get(field)
        if value is not None:
            point = point.field(field, float(value))
            field_count += 1

    if field_count == 0:
        logger.warning(f"Message for mote {mote_id} has no valid fields")
        return None

    # Use updated_timestamp (2025+ mapped data) for InfluxDB
    ts = message.get("updated_timestamp")
    if ts is None:
        # Fallback to regular timestamp if updated_timestamp not available
        ts = message.get("timestamp")

    if ts is not None:
        try:
            # Parse ISO format timestamp string to datetime
            if isinstance(ts, str):
                ts = datetime.fromisoformat(ts.replace('Z', '+00:00'))
            point = point.time(ts, WritePrecision.NS)
        except Exception as e:
            logger.error(f"Failed to parse timestamp '{ts}': {e}")
            return None

    return point

def write_batch(self, messages: Iterable[Dict]) -> None:
    """Write a batch of messages to InfluxDB 3.x.

    Args:

```

```

        messages: Iterable of sensor reading dictionaries
    """
    points: List[Point] = []
    skipped = 0

    for message in messages:
        point = self._to_point(message)
        if point is not None:
            points.append(point)
        else:
            skipped += 1

    if not points:
        logger.warning(f"No valid points to write (skipped {skipped}
invalid messages)")
        return

    try:
        logger.info(f"Writing {len(points)} points to InfluxDB
bucket={self.bucket}, org={self.org}, measurement={self.measurement}")
        if skipped > 0:
            logger.warning(f"Skipped {skipped} invalid messages")

        # InfluxDB 2.x write API
        self.write_api.write(bucket=self.bucket, org=self.org,
record=points)
        logger.info(f"Successfully wrote {len(points)} points to
measurement '{self.measurement}'")
    except Exception as e:
        logger.exception(f"Failed to write batch to InfluxDB: {e}")
        # Print first point for debugging
        if points:
            logger.error(f"Sample point that failed:
{points[0].to_line_protocol()}")
            raise # Re-raise to indicate failure

```

```

def close(self) -> None:
    """Close InfluxDB client connection."""
    try:
        self.client.close()
        logger.info("InfluxDB client closed")
    except Exception:
        logger.exception("Failed to close InfluxDB client")

```

[src/app/consumer/kafka_consumer.py](#)

"""Kafka consumer that batches messages for downstream processing."""

import json


```

import logging
import time
from typing import Dict, Iterable, List, Optional

from kafka import KafkaConsumer

logger = logging.getLogger(__name__)

class KafkaBatchConsumer:
    """Consume Kafka messages and yield them in batches."""

    def __init__(
        self,
        broker: str,
        topic: str,
        group_id: str = "sieis-consumers",
        batch_size: int = 100,
        batch_timeout: float = 1.0,
        auto_offset_reset: str = "earliest", # "latest" caused messages to be
missed if consumer starts after simulator
    ) -> None:
        self.topic = topic
        self.batch_size = batch_size
        self.batch_timeout = batch_timeout
        self.consumer = KafkaConsumer(
            topic,
            bootstrap_servers=[broker],
            group_id=group_id,
            auto_offset_reset=auto_offset_reset,
            enable_auto_commit=False,
            value_deserializer=lambda m: json.loads(m.decode("utf-8")),
        )

    def consume_batches(self) -> Iterable[List[Dict]]:
        """Yield message batches based on size or time window."""
        batch: List[Dict] = []
        batch_start = time.time()

        while True:
            records = self.consumer.poll(timeout_ms=100)
            for _, msgs in records.items():
                for msg in msgs:
                    if msg.value is None:
                        continue
                    batch.append(msg.value)

            now = time.time()
            if batch and (len(batch) >= self.batch_size or (now - batch_start)
>= self.batch_timeout):

```

```
        yield batch
        batch = []
        batch_start = now

def commit(self) -> None:
    """Commit current offsets after successful processing."""
    try:
        self.consumer.commit()
    except Exception:
        logger.exception("Failed to commit Kafka offsets")

def close(self) -> None:
    try:
        self.consumer.close()
    except Exception:
        logger.exception("Failed to close Kafka consumer")
```

src/app/consumer/main.py

```
"""Kafka-to-InfluxDB-and-MinIO consumer entry point with dual-write
pattern."""

import logging
from minio import Minio

from src.app.config import (
    INFLUX_BUCKET,
    INFLUX_ORG,
    INFLUX_TOKEN,
    INFLUX_URL,
    KAFKA_BROKER,
    KAFKA_TOPIC,
    MINIO_ACCESS_KEY,
    MINIO_BUCKET,
    MINIO_ENDPOINT,
    MINIO_SECRET_KEY,
    MINIO_SECURE,
)
from src.app.consumer.kafka_consumer import KafkaBatchConsumer
from src.app.consumer.influx_writer import InfluxWriter
from src.app.consumer.parquet_writer import ParquetWriter

logger = logging.getLogger(__name__)

def run_consumer() -> None:
    """Run dual-write consumer: Kafka → InfluxDB (hot, real-time) + MinIO
    (cold, Parquet)."""
    logging.basicConfig(level=logging.INFO)

    # Initialize Kafka consumer
    consumer = KafkaBatchConsumer(broker=KAFKA_BROKER, topic=KAFKA_TOPIC)

    # Initialize InfluxDB 2.x writer (hot path - real-time queries)
    influx_writer = InfluxWriter(
        url=INFLUX_URL,
        token=INFLUX_TOKEN,
        org=INFLUX_ORG,
        bucket=INFLUX_BUCKET,
    )

    # Initialize MinIO client
    minio_client = Minio(
        endpoint=MINIO_ENDPOINT,
        access_key=MINIO_ACCESS_KEY,
        secret_key=MINIO_SECRET_KEY,
        secure=MINIO_SECURE
```

```

)

# Initialize Parquet writer (cold path - historical archive)
parquet_writer = ParquetWriter(
    minio_client=minio_client,
    bucket_name=MINIO_BUCKET
)

logger.info("Consumer started with dual-write pattern: InfluxDB + MinIO")
logger.info(f"  InfluxDB: {INFLUX_URL} (org={INFLUX_ORG},
bucket={INFLUX_BUCKET})")
logger.info(f"  MinIO: {MINIO_ENDPOINT} (bucket={MINIO_BUCKET})")

try:
    for batch in consumer.consume_batches():
        # Dual-write pattern: write to both destinations
        influx_success = False
        parquet_success = False

        # Write to InfluxDB (hot path - critical for real-time dashboard)
        try:
            influx_writer.write_batch(batch)
            influx_success = True
        except Exception as e:
            logger.error(f"InfluxDB write failed: {e}")
            # Continue to try MinIO write even if InfluxDB fails

        # Write to MinIO (cold path - best effort, don't block on failure)
        try:
            parquet_writer.write_batch(batch)
            parquet_success = True
        except Exception as e:
            logger.warning(f"MinIO write failed (non-critical): {e}")
            # Don't block consumer on MinIO failures

        # Commit offset only if InfluxDB write succeeded (critical path)
        if influx_success:
            consumer.commit()
            if parquet_success:
                logger.debug("Dual-write successful: InfluxDB + MinIO")
            else:
                logger.warning("Partial write: InfluxDB succeeded, MinIO
failed")
        else:
            logger.error("Skipping Kafka commit due to InfluxDB write
failure")

    except KeyboardInterrupt:
        logger.info("KeyboardInterrupt received - stopping consumer")
    finally:

```

```

        consumer.close()
        influx_writer.close()
        parquet_writer.close()
        logger.info("Consumer shutdown complete")

if __name__ == "__main__":
    run_consumer()

```

src/app/consumer/parquet_writer.py

```

"""MinIO Parquet writer for historical sensor data archive."""

```

```

import logging
import hashlib
from io import BytesIO
from typing import Dict, List, Set, Tuple
from datetime import datetime, timezone

```

```

import pandas as pd
import pyarrow as pa
import pyarrow.parquet as pq
from minio import Minio
from minio.error import S3Error

```

```

logger = logging.getLogger(__name__)

```

```

class ParquetWriter:

```

```

    """Write sensor data batches to MinIO as Parquet files with date
    partitioning."""

```

```

    def __init__(self, minio_client: Minio, bucket_name: str,
        enable_deduplication: bool = True) -> None:

```

```

        """Initialize ParquetWriter with MinIO client.

```

```

        Args:

```

```

            minio_client: Initialized MinIO client

```

```

            bucket_name: Target bucket name (e.g., 'sieis-archive')

```

```

            enable_deduplication: Enable in-memory deduplication of messages

```

```

        (default: True)

```

```

        """

```

```

        self.minio_client = minio_client

```

```

        self.bucket_name = bucket_name

```

```

        self.enable_deduplication = enable_deduplication

```

```

        # In-memory cache of seen (mote_id, timestamp) pairs for deduplication

```

```

        self._seen_keys: Set[Tuple[int, str]] = set()

```

```

# Ensure bucket exists
if not self.minio_client.bucket_exists(bucket_name):
    logger.warning(f"Bucket '{bucket_name}' does not exist, attempting
to create")
    try:
        self.minio_client.make_bucket(bucket_name)
        logger.info(f"Created bucket '{bucket_name}'")
    except Exception as e:
        logger.error(f"Failed to create bucket '{bucket_name}': {e}")
        raise

```

```

def _deduplicate_messages(self, messages: List[Dict]) -> List[Dict]:
    """Remove duplicate messages based on (mote_id, timestamp) key.

```

```

    Args:
        messages: List of sensor reading dictionaries

```

```

    Returns:
        Deduplicated list of messages
    """

```

```

    if not self.enable_deduplication:
        return messages

```

```

    deduplicated = []
    initial_count = len(messages)

```

```

    for msg in messages:
        # Use updated_timestamp if available, otherwise regular timestamp
        ts = msg.get("updated_timestamp") or msg.get("timestamp")
        if ts is None:
            logger.warning(f"Message missing timestamp, including anyway:
{msg}")
            deduplicated.append(msg)
            continue

```

```

        key = (msg.get("mote_id"), str(ts))
        if key not in self._seen_keys:
            self._seen_keys.add(key)
            deduplicated.append(msg)

```

```

    if initial_count > len(deduplicated):
        logger.info(f"Deduplicated {initial_count - len(deduplicated)}
duplicate messages")

```

```

    return deduplicated

```

```

def clear_deduplication_cache(self) -> None:
    """Clear the in-memory deduplication cache.

```

Call this periodically to prevent unbounded memory growth in long-

```

running processes.
"""
    logger.info(f"Clearing deduplication cache ({len(self._seen_keys)}
entries)")
    self._seen_keys.clear()

def write_batch(self, messages: List[Dict]) -> None:
    """Convert message batch to Parquet and upload to MinIO.

    Messages are partitioned by date extracted from updated_timestamp.
    File path structure:
    year=YYYY/month=MM/day=DD/mote_id=X/batch_TIMESTAMP.parquet

    Uses deterministic filenames based on content to ensure idempotency.
    Deduplicates messages to prevent duplicate data.

    Args:
        messages: List of sensor reading dictionaries
    """
    if not messages:
        logger.warning("Empty batch provided to ParquetWriter, skipping")
        return

    try:
        # Deduplicate messages first
        messages = self._deduplicate_messages(messages)

        if not messages:
            logger.info("All messages were duplicates, skipping write")
            return

        # Convert to DataFrame
        df = pd.DataFrame(messages)

        # Parse updated_timestamp for partitioning
        if "updated_timestamp" not in df.columns:
            logger.error("Messages missing 'updated_timestamp' field,
cannot partition")
            return

        df["updated_timestamp"] = pd.to_datetime(df["updated_timestamp"])

        # Extract partition columns
        df["year"] = df["updated_timestamp"].dt.year
        df["month"] =
df["updated_timestamp"].dt.month.astype(str).str.zfill(2)
        df["day"] =
df["updated_timestamp"].dt.day.astype(str).str.zfill(2)

        # Group by partition and mote_id for efficient storage

```

```

grouped = df.groupby(["year", "month", "day", "mote_id"])

upload_count = 0
skipped_count = 0
for (year, month, day, mote_id), group in grouped:
    try:
        # Create deterministic filename based on content
        (idempotent)
        min_ts =
group["updated_timestamp"].min().strftime("%Y%m%d_%H%M%S")
        max_ts =
group["updated_timestamp"].max().strftime("%Y%m%d_%H%M%S")
        record_count = len(group)

        # Create content hash for additional uniqueness (in case
        same time range)
        content_str =
f"{min_ts}_{max_ts}_{record_count}_{mote_id}"
        content_hash =
hashlib.sha256(content_str.encode()).hexdigest()[:8]

        object_name =
f"year={year}/month={month}/day={day}/mote_id={mote_id}/batch_{min_ts}_{max_ts}
_{record_count}_{content_hash}.parquet"

        # Check if file already exists (idempotency - skip if
        exists)
        try:
            self.minio_client.stat_object(self.bucket_name,
object_name)
            logger.debug(f"File {object_name} already exists,
skipping upload (idempotent)")
            skipped_count += 1
            continue
        except S3Error as e:
            if e.code != "NoSuchKey":
                logger.error(f"Error checking object existence for
{object_name}: {e}")
                # If checking fails, try to proceed with upload
                anyway or continue?
                # Safer to log and attempt upload, or just raise.
                # Let's attempt upload, worst case it overwrites
                or fails.

                # File doesn't exist, proceed with upload

        # Drop partition columns (already in path) and write to
        buffer

        parquet_df = group.drop(columns=["year", "month", "day"])
        buffer = BytesIO()

```



```

        # Write Parquet to in-memory buffer
        table = pa.Table.from_pandas(parquet_df)
        pq.write_table(table, buffer, compression="snappy")

        # Upload to MinIO
        buffer.seek(0)
        self.minio_client.put_object(
            bucket_name=self.bucket_name,
            object_name=object_name,
            data=buffer,
            length=buffer.getbuffer().nbytes,
            content_type="application/octet-stream"
        )

        upload_count += 1
        logger.debug(f"Uploaded {len(group)} records to
{object_name}")

    except Exception as e:
        logger.error(f"Failed to process group for mote {mote_id}
(year={year}, month={month}, day={day}): {e}")
        # Continue to next group instead of failing the whole
batch
        continue

    if skipped_count > 0:
        logger.info(f"Successfully wrote {len(messages)} messages to
MinIO: {upload_count} new files, {skipped_count} skipped (already exist)")
    else:
        logger.info(f"Successfully wrote {len(messages)} messages to
MinIO as {upload_count} Parquet files")

    except Exception as e:
        logger.exception(f"Failed to write batch to MinIO: {e}")
        # Don't raise - we don't want Parquet write failures to block
InfluxDB writes

    def close(self) -> None:
        """Cleanup resources (MinIO client doesn't need explicit closing)."""
        pass

```

[src/app/dashboard/__init__.py](#)

SIEIS Dashboard package

[src/app/dashboard/app.py](#)

"""SIEIS Streamlit Dashboard - main entry point.

Usage:

```
streamlit run src/app/dashboard/app.py
```

Or via the run script:

```
python scripts/run_dashboard.py
```

```
"""
```

```
import os
import sys
```

```
# Ensure project root is in path
sys.path.insert(0,
os.path.dirname(os.path.dirname(os.path.dirname(os.path.abspath(
__file__))))))
```

```
from dotenv import load_dotenv
load_dotenv()
```

```
import streamlit as st
```

```
st.set_page_config(
    page_title="SIEIS - Smart IoT Sensor Dashboard",
    page_icon="🌱",
    layout="wide",
    initial_sidebar_state="expanded",
)
```

```
# — Home page content
```

```
st.title("🌱 SIEIS - Smart IoT Environmental Information System")
st.markdown("---")
```

```
col1, col2, col3, col4 = st.columns(4)
```

```
with col1:
    st.metric("Data Source", "Intel Lab Dataset", delta="54 notes")
```

```
with col2:
    st.metric("Sensors", "4 metrics", delta="Temp / Humidity / Light /
Voltage")
```

```
with col3:
    st.metric("Storage", "Dual-write", delta="InfluxDB + MinIO")
```




```
with col4:
    st.metric("ML Model", "Isolation Forest", delta="Anomaly Detection")
```

```
st.markdown("---")
```

```

st.markdown("""
## Welcome to SIEIS Dashboard

Use the **sidebar** to navigate between views:

Page	Purpose
 Real-time Monitor	Live sensor data from InfluxDB (last 1h-24h)
 Historical Analysis	Long-term trends from MinIO Parquet archives
 Anomaly Detection	ML-powered anomaly analysis and alerts

### Quick Start
1. Make sure Docker containers are running: `docker compose up -d`
2. Verify data is flowing: `python scripts/verify_influxDb.py`
3. Train the ML model: `python scripts/train_model.py`
4. Start the API: `python -m src.app.api_server`

### System Architecture
```
Sensors → Simulator → Redpanda/Kafka → Consumer ↗ InfluxDB (hot) → API →
Dashboard
 ↳ MinIO (cold) → ML →
Dashboard
```
""")

st.sidebar.success("Select a page above 🙌")

```

src/app/dashboard/pages/1_Realtime_Monitor.py

```

"""Real-time sensor monitoring page – queries InfluxDB directly."""

import os
import sys
sys.path.insert(0,
os.path.dirname(os.path.dirname(os.path.dirname(os.path.dirname(
os.path.abspath(__file__))))))

from dotenv import load_dotenv
load_dotenv()

import time
import logging
from typing import List, Optional
import pandas as pd
import plotly.express as px
import plotly.graph_objects as go
import streamlit as st

```

```

logger = logging.getLogger(__name__)

st.set_page_config(page_title="Real-time Monitor", page_icon="🔴",
layout="wide")
st.title("🔴 Real-time Sensor Monitor")
st.caption("Live data from InfluxDB – auto-refreshes every 30 seconds")

from src.app import config

@st.cache_resource
def get_influx_client():
    from influxdb_client import InfluxDBClient
    return InfluxDBClient(url=config.INFLUX_URL, token=config.INFLUX_TOKEN,
org=config.INFLUX_ORG)

def query_influx(flux: str) -> pd.DataFrame:
    """Run a Flux query and return results as DataFrame."""
    try:
        client = get_influx_client()
        query_api = client.query_api()
        tables = query_api.query_data_frame(flux, org=config.INFLUX_ORG)
        if isinstance(tables, list):
            if not tables:
                return pd.DataFrame()
            df = pd.concat(tables, ignore_index=True)
        else:
            df = tables
        # Drop internal InfluxDB columns
        drop_cols = [c for c in df.columns if c.startswith("_") and c not in
["_value", "_time", "_field"]]
        df = df.drop(columns=drop_cols, errors="ignore")
        return df
    except Exception as e:
        st.warning(f"InfluxDB query failed: {e}")
        return pd.DataFrame()

# — Sidebar Controls


---


with st.sidebar:
    st.header("⚙️ Controls")
    time_window = st.selectbox(
        "Time window",
        options=["-15m", "-1h", "-6h", "-24h"],
        index=1,
        format_func=lambda x: {"-15m": "Last 15 min", "-1h": "Last 1 hour",
"-6h": "Last 6 hours", "-24h": "Last 24 hours"}.get(x,

```

```

x),
    )
    metric = st.selectbox("Metric to plot", ["temperature", "humidity",
"light", "voltage"])
    auto_refresh = st.checkbox("Auto-refresh (30s)", value=False)

    st.markdown("---")
    st.markdown("***Connection**")
    st.code(f"InfluxDB: {config.INFLUX_URL}\nBucket: {config.INFLUX_BUCKET}",
language="text")

# — KPI Row


---



st.subheader("📊 Key Metrics")

flux_active = f"""
from(bucket: "{config.INFLUX_BUCKET}")
  |> range(start: -15m)
  |> filter(fn: (r) => r["_measurement"] == "sensor_reading")
  |> keep(columns: ["mote_id"])
  |> distinct(column: "mote_id")
  |> count()
"""

flux_summary = f"""
from(bucket: "{config.INFLUX_BUCKET}")
  |> range(start: {time_window})
  |> filter(fn: (r) => r["_measurement"] == "sensor_reading")
  |> filter(fn: (r) => r["_field"] == "{metric}")
  |> mean()
"""

col1, col2, col3, col4 = st.columns(4)

df_active = query_influx(flux_active)
active_count = 0
if not df_active.empty and "_value" in df_active.columns:
    active_count = int(df_active["_value"].iloc[0])
elif not df_active.empty:
    active_count = len(df_active)

col1.metric("Active Motes (15m)", active_count, delta=None)

df_avg = query_influx(flux_summary)
if not df_avg.empty and "_value" in df_avg.columns:
    avg_val = round(df_avg["_value"].mean(), 2)
    col2.metric(f"Avg {metric.title()}", avg_val)
else:
    col2.metric(f"Avg {metric.title()}", "N/A")

```

```


# Count total records
flux_count = f"""
from(bucket: "{config.INFLUX_BUCKET}")
  |> range(start: {time_window})
  |> filter(fn: (r) => r["_measurement"] == "sensor_reading")
  |> filter(fn: (r) => r["_field"] == "temperature")
  |> count()
"""

df_count = query_influx(flux_count)
total_recs = 0
if not df_count.empty and "_value" in df_count.columns:
    total_recs = int(df_count["_value"].sum())
col3.metric("Total Records", f"{total_recs:,}")
col4.metric("Time Window", time_window)

# — Time Series Chart


---



st.subheader(f" {metric.title()} Over Time")

flux_ts = f"""
from(bucket: "{config.INFLUX_BUCKET}")
  |> range(start: {time_window})
  |> filter(fn: (r) => r["_measurement"] == "sensor_reading")
  |> filter(fn: (r) => r["_field"] == "{metric}")
  |> aggregateWindow(every: 1m, fn: mean, createEmpty: false)
"""

df_ts = query_influx(flux_ts)

if not df_ts.empty and "_value" in df_ts.columns:
    time_col = "_time" if "_time" in df_ts.columns else df_ts.columns[0]
    mote_col = "mote_id" if "mote_id" in df_ts.columns else None

    if mote_col:
        # Show top 5 motes to avoid clutter
        top_motes =
df_ts.groupby(mote_col)["_value"].count().nlargest(5).index.tolist()
        df_plot = df_ts[df_ts[mote_col].isin(top_motes)]
        fig = px.line(
            df_plot, x=time_col, y="_value", color=mote_col,
            title=f"{metric.title()} – Top 5 most active motes",
            labels={"_value": metric, time_col: "Time", mote_col: "Mote"},
        )
    else:
        fig = px.line(df_ts, x=time_col, y="_value",
            title=f"{metric.title()} over time",
            labels={"_value": metric})

```

```

        fig.update_layout(height=400, legend=dict(orientation="h",
yanchor="bottom", y=1.02))
        st.plotly_chart(fig, use_container_width=True)
    else:
        st.info("No data available for the selected time window. Make sure the
simulator and consumer are running.")

# — Latest readings grid


---


st.subheader("📊 Latest Reading per Mote")

flux_latest = f"""
from(bucket: "{config.INFLUX_BUCKET}")
  |> range(start: {time_window})
  |> filter(fn: (r) => r["_measurement"] == "sensor_reading")
  |> last()
  |> pivot(rowKey:["_time","mote_id"], columnKey: ["_field"], valueColumn:
"_value")
"""

df_latest = query_influx(flux_latest)
if not df_latest.empty:
    display_cols = [c for c in ["mote_id", "temperature", "humidity", "light",
"voltage", "_time"] if c in df_latest.columns]
    df_display = df_latest[display_cols].copy()
    if "_time" in df_display.columns:
        df_display = df_display.rename(columns={"_time": "last_seen"})
    for col in ["temperature", "humidity", "light", "voltage"]:
        if col in df_display.columns:
            df_display[col] = df_display[col].round(2)
    st.dataframe(df_display.sort_values("mote_id") if "mote_id" in
df_display.columns else df_display, use_container_width=True)
else:
    st.info("No recent data. Ensure the data pipeline is running.")

# — Auto-refresh


---


if auto_refresh:
    time.sleep(30)
    st.rerun()

```

[src/app/dashboard/pages/2_Historical_Analysis.py](#)

"""Historical data analysis – queries MinIO Parquet files."""

```

import os
import sys
sys.path.insert(0,
os.path.dirname(os.path.dirname(os.path.dirname(os.path.dirname(

```

```

e(os.path.abspath(__file__))))))

from dotenv import load_dotenv
load_dotenv()

import io
import logging
import pandas as pd
import plotly.express as px
import plotly.graph_objects as go
import streamlit as st

logger = logging.getLogger(__name__)

st.set_page_config(page_title="Historical Analysis", page_icon="📊",
layout="wide")
st.title("📊 Historical Sensor Analysis")
st.caption("Long-term trends from MinIO cold storage (Parquet files)")

from src.app import config

@st.cache_data(ttl=300)
def load_local_data(max_rows: int = 100_000) -> pd.DataFrame:
    """Load from local historical_data.txt (fallback when MinIO is empty)."""
    hist_path = config.DATA_DIR / "processed" / "historical_data.txt"
    if not hist_path.exists():
        return pd.DataFrame()
    cols = ["date", "time", "epoch", "mote_id", "temperature", "humidity",
"light", "voltage", "updated_timestamp"]
    df = pd.read_csv(str(hist_path), sep=r"\s+", names=cols, na_values=["N/A",
"nan", ""], nrows=max_rows, on_bad_lines="skip")
    df["timestamp"] = pd.to_datetime(df["date"] + " " + df["time"],
errors="coerce")
    # Normalize mote_id to plain int - mixed int/float parsing causes
duplicates (e.g. 1 vs 1.0)
    df["mote_id"] = pd.to_numeric(df["mote_id"],
errors="coerce").astype("Int64")
    for col in ["temperature", "humidity", "light", "voltage"]:
        df[col] = pd.to_numeric(df[col], errors="coerce")
    df = df.dropna(subset=["temperature", "humidity"], how="all")
    return df

@st.cache_data(ttl=300)
def load_minio_data(max_rows: int = 500_000) -> pd.DataFrame:
    """Load ALL Parquet files from MinIO in parallel, capped at max_rows total
rows."""
    from concurrent.futures import ThreadPoolExecutor, as_completed

```



```

from minio import Minio

try:
    client = Minio(
        config.MINIO_ENDPOINT,
        access_key=config.MINIO_ACCESS_KEY,
        secret_key=config.MINIO_SECRET_KEY,
        secure=config.MINIO_SECURE,
    )
    objects = list(client.list_objects(config.MINIO_BUCKET,
recursive=True))
    parquet_objects = sorted(
        [o for o in objects if o.object_name.endswith(".parquet")],
        key=lambda o: o.object_name,
    )
    if not parquet_objects:
        return pd.DataFrame()

    def _fetch(obj_name: str):
        """Download one Parquet file and return a DataFrame."""
        try:
            # Each thread needs its own client instance (not thread-safe
to share)

            _client = Minio(
                config.MINIO_ENDPOINT,
                access_key=config.MINIO_ACCESS_KEY,
                secret_key=config.MINIO_SECRET_KEY,
                secure=config.MINIO_SECURE,
            )
            resp = _client.get_object(config.MINIO_BUCKET, obj_name)
            data = resp.read()
            resp.close()
            return pd.read_parquet(io.BytesIO(data))
        except Exception as e:
            logger.warning(f"Failed to load {obj_name}: {e}")
            return pd.DataFrame()

    frames = []
    total_rows = 0
    # Use up to 8 parallel workers - saturates MinIO I/O without
overwhelming it
    with ThreadPoolExecutor(max_workers=8) as pool:
        futures = {
            pool.submit(_fetch, obj.object_name): obj.object_name
            for obj in parquet_objects
        }
        for future in as_completed(futures):
            if total_rows >= max_rows:
                future.cancel()
                continue

```

```

        chunk = future.result()
        if not chunk.empty:
            frames.append(chunk)
            total_rows += len(chunk)

    if not frames:
        return pd.DataFrame()

    df = pd.concat(frames, ignore_index=True)
    # Sort chronologically after parallel assembly
    if "timestamp" in df.columns:
        df = df.sort_values("timestamp").reset_index(drop=True)
    # Apply row cap
    if len(df) > max_rows:
        df = df.iloc[:max_rows]
    # Normalize mote_id to plain int
    if "mote_id" in df.columns:
        df["mote_id"] = pd.to_numeric(df["mote_id"],
errors="coerce").astype("Int64")
    return df

except Exception as e:
    logger.warning(f"MinIO load failed: {e}")
    return pd.DataFrame()

# — Sidebar Controls


---


import datetime as _dt

# Retrieve saved data bounds from session_state (populated after first load)
_saved_min: _dt.date | None = st.session_state.get("data_min")
_saved_max: _dt.date | None = st.session_state.get("data_max")

with st.sidebar:
    st.header("⚙ Controls")
    data_source = st.radio("Data source", ["Local file", "MinIO Parquet"])
    max_rows = st.slider("Max rows to load", 10_000, 2_000_000, 500_000,
10_000)
    selected_metric = st.selectbox("Metric", ["temperature", "humidity",
"light", "voltage"])

    st.markdown("---")
    st.markdown("***📅 Date Range Filter***")
    if _saved_min and _saved_max:
        date_start = st.date_input(
            "Start Date", value=_saved_min,
            min_value=_saved_min, max_value=_saved_max,
            key="date_start",

```

```

    )
    date_end = st.date_input(
        "End Date", value=_saved_max,
        min_value=_saved_min, max_value=_saved_max,
        key="date_end",
    )
else:
    st.caption("Loading data to determine date range...")
    date_start = date_end = None

    st.markdown("---")
    if data_source == "MinIO Parquet":
        st.code(f"MinIO: {config.MINIO_ENDPOINT}\nBucket:
{config.MINIO_BUCKET}", language="text")
        st.caption("All Parquet files loaded. Use the slider to limit total
rows.")
    else:
        st.code("File: data/processed/historical_data.txt", language="text")

# — Load data


---


_cache_col, _btn_col = st.columns([5, 1])
with _btn_col:
    if st.button("🔄 Refresh", help="Clear cache and reload data from
source"):
        load_minio_data.clear()
        load_local_data.clear()
        st.session_state.pop("data_min", None)
        st.session_state.pop("data_max", None)
        st.rerun()

if data_source == "MinIO Parquet":
    with st.spinner("Loading Parquet files from MinIO in parallel — results
are cached for 5 min after first load..."):
        df = load_minio_data(max_rows)
        if df.empty:
            st.warning("No Parquet files found in MinIO. Loading from local
file instead.")
            df = load_local_data(max_rows)
else:
    with st.spinner("Loading local data file..."):
        df = load_local_data(max_rows)

if df.empty:
    st.error("No data available. Run `python scripts/load_historical_data.py`
first.")
    st.stop()

# Ensure timestamp column

```

```

if "timestamp" not in df.columns and "date" in df.columns and "time" in
df.columns:
    df["timestamp"] = pd.to_datetime(df["date"] + " " + df["time"],
errors="coerce")
elif "timestamp" not in df.columns and "_time" in df.columns:
    df["timestamp"] = pd.to_datetime(df["_time"], errors="coerce")

# Force timestamp to proper datetime dtype regardless of source
if "timestamp" in df.columns:
    df["timestamp"] = pd.to_datetime(df["timestamp"], errors="coerce")

# Ensure numeric sensor cols
for col in ["temperature", "humidity", "light", "voltage"]:
    if col in df.columns:
        df[col] = pd.to_numeric(df[col], errors="coerce")

# — Save data bounds to session_state, rerun once so pickers appear


---


if "timestamp" in df.columns:
    _ts = df["timestamp"].dropna()
    if len(_ts) > 0:
        _new_min, _new_max = _ts.min().date(), _ts.max().date()
        if st.session_state.get("data_min") != _new_min or
st.session_state.get("data_max") != _new_max:
            st.session_state["data_min"] = _new_min
            st.session_state["data_max"] = _new_max
            st.rerun() # re-render sidebar with correct picker bounds

# — Apply date filter


---


if date_start and date_end and "timestamp" in df.columns:
    mask = (df["timestamp"].dt.date >= date_start) & (df["timestamp"].dt.date
<= date_end)
    df = df[mask].reset_index(drop=True)
    if df.empty:
        st.warning("No records in the selected date range. Adjust the filter
in the sidebar.")
        st.stop()

# — Sidebar stats (injected after filter applied)


---


with st.sidebar:
    st.markdown("***🇮🇹 Dataset Stats**")
    st.metric("Records", f"{len(df):,}")
    st.metric("Unique Motes", df["mote_id"].nunique() if "mote_id" in
df.columns else "N/A")

# — KPI Row


---



```

```

st.subheader("📊 Dataset Overview")
col1, col2, col3, col4 = st.columns(4)
col1.metric("Total Records", f"{len(df):,}")
col2.metric("Unique Motes", df["mote_id"].nunique() if "mote_id" in df.columns
else "N/A")
if "timestamp" in df.columns:
    ts = df["timestamp"].dropna()
    if len(ts) > 0:
        col3.metric("Date Range Start", str(ts.min().date()))
        col4.metric("Date Range End", str(ts.max().date()))

```

— Daily averages chart

```

st.subheader(f"📈 Daily Average {selected_metric.title()}")

if "timestamp" in df.columns and selected_metric in df.columns:
    df_daily = df.copy()
    df_daily["date"] = df_daily["timestamp"].dt.date
    daily_avg = df_daily.groupby("date")[selected_metric].mean().reset_index()
    daily_avg.columns = ["date", "avg_value"]
    daily_avg = daily_avg.dropna()

    if not daily_avg.empty:
        fig_daily = px.line(
            daily_avg, x="date", y="avg_value",
            title=f"Daily Average {selected_metric.title()}",
            labels={"avg_value": selected_metric, "date": "Date"},
        )
        fig_daily.update_traces(line_color="#1f77b4")
        fig_daily.update_layout(height=350)
        st.plotly_chart(fig_daily, use_container_width=True)
    else:
        st.info("No daily data to display.")
else:
    st.info("Timestamp or metric column not available in this dataset.")

```

— Hourly heatmap

```

st.subheader(f"🌡️ Hourly Seasonality Heatmap - {selected_metric.title()}")

if "timestamp" in df.columns and selected_metric in df.columns:
    df_heat = df.copy()
    df_heat["hour"] = df_heat["timestamp"].dt.hour
    df_heat["day_of_week"] = df_heat["timestamp"].dt.day_name()

    heat_pivot = (
        df_heat.groupby(["day_of_week", "hour"])[selected_metric]
        .mean()
        .reset_index()
    )

```

```


        .pivot(index="day_of_week", columns="hour", values=selected_metric)
    )

    day_order = ["Monday", "Tuesday", "Wednesday", "Thursday", "Friday",
"Saturday", "Sunday"]
    heat_pivot = heat_pivot.reindex([d for d in day_order if d in
heat_pivot.index])

    if not heat_pivot.empty:
        fig_heat = px.imshow(
            heat_pivot,
            title=f"{selected_metric.title()} by Hour & Day of Week",
            labels={"x": "Hour of Day", "y": "Day", "color": selected_metric},
            color_continuous_scale="RdYlGn",
            aspect="auto",
        )
        fig_heat.update_layout(height=350)
        st.plotly_chart(fig_heat, use_container_width=True)

# — Distribution


---


st.subheader(f" Value Distribution - {selected_metric.title()}")

if selected_metric in df.columns:
    col_a, col_b = st.columns(2)

    with col_a:
        fig_hist = px.histogram(
            df.dropna(subset=[selected_metric]),
            x=selected_metric, nbins=50,
            title=f"Distribution of {selected_metric}",
            color_discrete_sequence=["#636EFA"],
        )
        fig_hist.update_layout(height=350)
        st.plotly_chart(fig_hist, use_container_width=True)

    with col_b:
        if "mote_id" in df.columns:
            box_df = df.dropna(subset=[selected_metric])
            top_motes = box_df["mote_id"].value_counts().nlargest(8).index
            box_df = box_df[box_df["mote_id"].isin(top_motes)]
            fig_box = px.box(
                box_df, x="mote_id", y=selected_metric,
                title=f"{selected_metric.title()} by Mote (top 8)",
                color="mote_id",
            )
            fig_box.update_layout(height=350, showlegend=False)
            st.plotly_chart(fig_box, use_container_width=True)

```

```
# — Raw data preview
```

```
with st.expander("🔍 Raw data sample (first 200 rows)":  
    preview_cols = [c for c in ["mote_id", "timestamp", "temperature",  
"humidity", "light", "voltage"] if c in df.columns]  
    st.dataframe(df[preview_cols].head(200), use_container_width=True)
```

```
src/app/dashboard/pages/3_Anomaly_Detection.py
```

```
"""Anomaly detection page – ML model inference via the SIEIS API."""
```

```
import os  
import sys  
sys.path.insert(0,  
os.path.dirname(os.path.dirname(os.path.dirname(os.path.dirname  
os.path.abspath(__file__))))))
```

```
from dotenv import load_dotenv  
load_dotenv()
```

```
import json  
import logging  
import pandas as pd  
import plotly.express as px  
import plotly.graph_objects as go  
import requests  
import streamlit as st
```

```
logger = logging.getLogger(__name__)
```

```
st.set_page_config(page_title="Anomaly Detection", page_icon="🤖",  
layout="wide")  
st.title("🤖 Anomaly Detection")  
st.caption("ML-powered anomaly analysis using Isolation Forest")
```

```
from src.app import config
```

```
_api_url = os.getenv("API_URL", f"http://localhost:{config.API_PORT}")  
API_BASE = f"{_api_url}/api/v1"
```

```
def call_api(path: str, method: str = "GET", body: dict = None):  
    """Call the SIEIS API with error handling."""  
    try:  
        url = f"{API_BASE}{path}"  
        if method == "POST":  
            resp = requests.post(url, json=body, timeout=5)  
        else:  
            resp = requests.get(url, timeout=5)
```

```

        resp.raise_for_status()
        return resp.json(), None
    except requests.exceptions.ConnectionError:
        return None, "API not running. Start it with: python -m
src.app.api_server"
    except Exception as e:
        return None, str(e)

# — API Status


---


health_data, health_err = call_api("/health", method="GET")

col_status, col_model = st.columns(2)
with col_status:
    if health_err:
        st.error(f"❌ API Offline - {health_err}")
        api_online = False
    else:
        st.success("✅ API Online")
        api_online = True

with col_model:
    if api_online and health_data:
        model_loaded = health_data.get("model_loaded", False)
        if model_loaded:
            st.success("✅ ML Model Loaded")
        else:
            st.warning("⚠️ Model not trained yet - using rule-based fallback")
            st.caption("Run: `python scripts/train_model.py` to train the
model")

st.markdown("----")

# — Model Info


---


if api_online:
    with st.expander("📄 Model Info", expanded=False):
        model_info, err = call_api("/ml/model/info")
        if model_info:
            st.json(model_info)
        else:
            st.warning(err or "Could not fetch model info")

col_reload, _ = st.columns([1, 3])
with col_reload:
    if st.button("🔄 Reload Model"):
        result, err = call_api("/ml/model/reload", method="POST")
        if result:

```



```

        st.success("Model reloaded!")
    else:
        st.error(err)

st.markdown("---")

# — Single Reading Prediction


---


st.subheader("🔍 Single Reading Prediction")
st.markdown("Enter sensor values to check if they're anomalous.")

with st.form("predict_form"):
    col1, col2, col3, col4 = st.columns(4)
    with col1:
        temperature = st.number_input("Temperature (°C)", min_value=-50.0,
max_value=100.0, value=22.5, step=0.1)
    with col2:
        humidity = st.number_input("Humidity (%)", min_value=0.0,
max_value=100.0, value=55.0, step=0.1)
    with col3:
        light = st.number_input("Light (Lux)", min_value=0.0,
max_value=5000.0, value=300.0, step=1.0)
    with col4:
        voltage = st.number_input("Voltage (V)", min_value=0.0,
max_value=5.0, value=2.9, step=0.01)

    mote_id_input = st.text_input("Mote ID (optional)", value="",
placeholder="e.g. 1")
    submitted = st.form_submit_button("🔍 Check for Anomaly",
use_container_width=True)

if submitted:
    if not api_online:
        st.error("API is not running. Cannot make predictions.")
    else:
        payload = {
            "temperature": temperature,
            "humidity": humidity,
            "light": light,
            "voltage": voltage,
        }
        if mote_id_input.strip():
            payload["mote_id"] = mote_id_input.strip()

        with st.spinner("Running anomaly detection..."):
            result, err = call_api("/ml/predict/anomaly", method="POST",
body=payload)

            if result:

```

```

severity = result.get("severity", "normal")
score = result.get("anomaly_score", 0)
is_anomaly = result.get("is_anomaly", False)

color_map = {"normal": "🟢", "warning": "🟡", "critical": "🔴"}
emoji = color_map.get(severity, "⬜")

res_col1, res_col2, res_col3 = st.columns(3)
res_col1.metric("Anomaly Score", f"{score:.4f}", delta=None)
res_col2.metric("Status", f"{emoji} {severity.upper()}")
res_col3.metric("Is Anomaly", "YES ⚠️" if is_anomaly else "NO ✅")

# Score gauge
fig_gauge = go.Figure(go.Indicator(
    mode="gauge+number",
    value=score * 100,
    domain={"x": [0, 1], "y": [0, 1]},
    title={"text": "Anomaly Score (0=Normal, 100=Anomaly)"},
    gauge={
        "axis": {"range": [0, 100]},
        "bar": {"color": "darkred" if is_anomaly else
"darkgreen"},
        "steps": [
            {"range": [0, 50], "color": "lightgreen"},
            {"range": [50, 65], "color": "yellow"},
            {"range": [65, 100], "color": "lightcoral"},
        ],
        "threshold": {"line": {"color": "red", "width": 4},
"thickness": 0.75, "value": 65},
    },
))
fig_gauge.update_layout(height=300)
st.plotly_chart(fig_gauge, use_container_width=True)
else:
    st.error(f"Prediction failed: {err}")

st.markdown("---")

# — Batch anomaly scan

st.subheader("🇩🇪 Batch Anomaly Scan – Recent Sensor Data")
st.markdown("Scan latest readings from InfluxDB for anomalies.")

if st.button("🚀 Run Batch Scan", disabled=not api_online):
    # Fetch latest readings from InfluxDB
    try:
        from influxdb_client import InfluxDBClient
        client = InfluxDBClient(url=config.INFLUX_URL,
token=config.INFLUX_TOKEN, org=config.INFLUX_ORG)

```

```

        query_api = client.query_api()
        flux = f"""
from(bucket: "{config.INFLUX_BUCKET}")
  |> range(start: -1h)
  |> filter(fn: (r) => r["_measurement"] == "sensor_reading")
  |> last()
  |> pivot(rowKey:["_time","mote_id"], columnKey: ["_field"], valueColumn:
"_value")
"""

        tables = query_api.query_data_frame(flux, org=config.INFLUX_ORG)
        client.close()

        if isinstance(tables, list):
            df = pd.concat(tables, ignore_index=True) if tables else
pd.DataFrame()
        else:
            df = tables

        if df.empty:
            st.warning("No recent data in InfluxDB. Run the simulator first.")
        else:
            results = []
            progress = st.progress(0)
            for i, row in df.iterrows():
                payload = {
                    "temperature": float(row.get("temperature", 22)),
                    "humidity": float(row.get("humidity", 50)),
                    "light": float(row.get("light", 100)),
                    "voltage": float(row.get("voltage", 2.9)),
                    "mote_id": str(row.get("mote_id", "unknown")),
                }
                result, _ = call_api("/ml/predict/anomaly", method="POST",
body=payload)
                if result:
                    results.append({
                        "mote_id": payload["mote_id"],
                        "temperature": payload["temperature"],
                        "humidity": payload["humidity"],
                        "light": payload["light"],
                        "voltage": payload["voltage"],
                        "anomaly_score": result["anomaly_score"],
                        "is_anomaly": result["is_anomaly"],
                        "severity": result["severity"],
                    })
                progress.progress(min(1.0, (i + 1) / max(len(df), 1)))

        if results:
            df_results = pd.DataFrame(results)
            anomalies = df_results[df_results["is_anomaly"]]

```

```

a_col1, a_col2, a_col3 = st.columns(3)
a_col1.metric("Motes Scanned", len(df_results))
a_col2.metric("Anomalies Found", len(anomalies))
a_col3.metric("Anomaly Rate",
f"{len(anomalies)/max(len(df_results),1):.1%}")

# Anomaly score scatter
fig_scatter = px.scatter(
    df_results, x="mote_id", y="anomaly_score",
    color="severity",
    color_discrete_map={"normal": "green", "warning":
"orange", "critical": "red"},
    title="Anomaly Scores per Mote",
    size_max=15,
)
fig_scatter.add_hline(y=0.5, line_dash="dash",
line_color="orange", annotation_text="Warning threshold")
fig_scatter.add_hline(y=0.65, line_dash="dash",
line_color="red", annotation_text="Critical threshold")
st.plotly_chart(fig_scatter, use_container_width=True)

# — Full results table with colour-coded rows —————
st.subheader("📁 All Readings")

# Row-level background colours by severity
def _row_color(row):
    colors = {
        "normal": "background-color: #d4edda", # green
        "warning": "background-color: #fff3cd", # amber
        "critical": "background-color: #f8d7da", # red
    }
    style = colors.get(row.get("severity", "normal"), "")
    return [style] * len(row)

# Add a Status emoji column at the front for quick scanning
display_df = df_results.copy()
display_df.insert(0, "Status", display_df["severity"].map({
    "normal": "🟢 Normal",
    "warning": "🟡 Warning",
    "critical": "🔴 Critical",
}))
display_df = display_df.sort_values("anomaly_score",
ascending=False)

st.dataframe(
    display_df.style.apply(_row_color, axis=1),
    use_container_width=True,
)

```

```
# Keep a focused anomaly-only section below
if not anomalies.empty:
    st.subheader("🚨 Anomalous Readings Only")
    st.dataframe(
        anomalies.sort_values("anomaly_score",
ascending=False),
        use_container_width=True,
    )
except Exception as e:
    st.error(f"Batch scan failed: {e}")
```

`src/app/ml/__init__.py`

src/app/ml/detector.py

```
"""Isolation Forest anomaly detector for SIEIS sensor data."""
```

```
import json
import logging
import pickle
from datetime import datetime
from pathlib import Path
from typing import Dict, Optional, Tuple

import numpy as np
import pandas as pd
from sklearn.ensemble import IsolationForest
from sklearn.preprocessing import StandardScaler
from sklearn.pipeline import Pipeline
```

```
from src.app import config
```

```
logger = logging.getLogger(__name__)
```

```
REGISTRY_FILE = config.MODELS_DIR / "model_registry.json"
```

```
def _load_registry() -> Dict:
    if REGISTRY_FILE.exists():
        with open(REGISTRY_FILE) as f:
            return json.load(f)
    return {"models": [], "latest": None}
```

```
def _save_registry(registry: Dict):
    config.MODELS_DIR.mkdir(parents=True, exist_ok=True)
    with open(REGISTRY_FILE, "w") as f:
        json.dump(registry, f, indent=2)
```

```
def train_anomaly_detector(
    X: pd.DataFrame,
    contamination: float = 0.05,
    n_estimators: int = 100,
    random_state: int = 42,
) -> Tuple[Pipeline, Dict]:
    """Train an Isolation Forest anomaly detector.
```

Analogy: Imagine a forest of decision trees. Each tree randomly partitions data points. Anomalies (outliers) get isolated in fewer splits than normal points – like a rotten apple that's easy to spot.

Args:

```

        X: Feature matrix (temperature, humidity, light, voltage, hour,
day_of_week)
        contamination: Expected fraction of anomalies in data (5% default)
        n_estimators: Number of trees in the forest
        random_state: For reproducibility

```

Returns:

```

    pipeline: Trained sklearn Pipeline (scaler + model)
    metrics: Training statistics dict

```

```

"""

```

```

    logger.info(f"Training Isolation Forest: n_samples={len(X)},
contamination={contamination}")

```

```

    pipeline = Pipeline([
        ("scaler", StandardScaler()),
        ("model", IsolationForest(
            n_estimators=n_estimators,
            contamination=contamination,
            random_state=random_state,
            n_jobs=-1,
        )),
    ])

```

```

    pipeline.fit(X)

```

```

    # Compute training metrics
    predictions = pipeline.predict(X)
    scores = pipeline.decision_function(X)
    n_anomalies = int((predictions == -1).sum())
    anomaly_ratio = n_anomalies / len(X)

```

```

    metrics = {
        "n_samples": len(X),
        "n_features": X.shape[1],
        "feature_names": list(X.columns),
        "n_anomalies_detected": n_anomalies,
        "anomaly_ratio": round(anomaly_ratio, 4),
        "contamination": contamination,
        "n_estimators": n_estimators,
        "score_mean": round(float(np.mean(scores)), 4),
        "score_std": round(float(np.std(scores)), 4),
        "trained_at": datetime.utcnow().isoformat(),
    }

```

```

    logger.info(f"Training complete: {n_anomalies}/{len(X)} anomalies
({anomaly_ratio:.1%})")
    return pipeline, metrics

```

```

def save_model(pipeline: Pipeline, metrics: Dict, tag: Optional[str] = None) -

```

```

> str:
    """Save trained model and update model registry.

    Returns the filename of the saved model.
    """
    config.MODELS_DIR.mkdir(parents=True, exist_ok=True)

    timestamp = datetime.utcnow().strftime("%Y%m%d_%H%M%S")
    tag_suffix = f"_{tag}" if tag else ""
    filename = f"anomaly_detector_{timestamp}{tag_suffix}.pkl"
    filepath = config.MODELS_DIR / filename

    artifact = {
        "pipeline": pipeline,
        "metrics": metrics,
        "version": timestamp,
    }

    with open(filepath, "wb") as f:
        pickle.dump(artifact, f)

    logger.info(f"Model saved: {filepath}")

    # Update registry
    registry = _load_registry()
    registry["models"].append({
        "filename": filename,
        "trained_at": metrics["trained_at"],
        "n_samples": metrics["n_samples"],
        "anomaly_ratio": metrics["anomaly_ratio"],
        "features": metrics["feature_names"],
    })
    registry["latest"] = filename
    _save_registry(registry)

    logger.info(f"Registry updated: latest={filename}")
    return filename


def load_latest_model() -> Optional[Pipeline]:
    """Load the latest trained model from registry."""
    registry = _load_registry()
    latest = registry.get("latest")
    if not latest:
        logger.warning("No model in registry")
        return None

    model_path = config.MODELS_DIR / latest
    if not model_path.exists():
        logger.error(f"Model file not found: {model_path}")

```



```

        return None

    with open(model_path, "rb") as f:
        artifact = pickle.load(f)

    logger.info(f"Loaded model: {latest}")
    return artifact["pipeline"]

```

`src/app/ml/models/model_registry.json`

```

{
  "models": [
    {
      "filename": "anomaly_detector_20260225_041931.pkl",
      "trained_at": "2026-02-25T04:19:31.445233",
      "n_samples": 295037,
      "anomaly_ratio": 0.01,
      "features": [
        "temperature",
        "humidity",
        "light",
        "voltage",
        "hour",
        "day_of_week"
      ]
    },
    {
      "filename": "anomaly_detector_20260225_053331.pkl",
      "trained_at": "2026-02-25T05:33:31.290149",
      "n_samples": 295037,
      "anomaly_ratio": 0.01,
      "features": [
        "temperature",
        "humidity",
        "light",
        "voltage",
        "hour",
        "day_of_week"
      ]
    },
    {
      "filename": "anomaly_detector_20260225_054911.pkl",
      "trained_at": "2026-02-25T05:49:10.736852",
      "n_samples": 499638,
      "anomaly_ratio": 0.05,
      "features": [
        "temperature",
        "humidity",
        "light",

```

```

        "voltage",
        "hour",
        "day_of_week"
    ]
},
{
    "filename": "anomaly_detector_20260225_061932.pkl",
    "trained_at": "2026-02-25T06:19:32.665523",
    "n_samples": 1740246,
    "anomaly_ratio": 0.05,
    "features": [
        "temperature",
        "humidity",
        "light",
        "voltage",
        "hour",
        "day_of_week"
    ]
},
{
    "filename": "anomaly_detector_20260225_104131_scheduled.pkl",
    "trained_at": "2026-02-25T10:41:31.463276",
    "n_samples": 1815412,
    "anomaly_ratio": 0.05,
    "features": [
        "temperature",
        "humidity",
        "light",
        "voltage",
        "hour",
        "day_of_week"
    ]
}
],
"latest": "anomaly_detector_20260225_104131_scheduled.pkl"
}

```

[src/app/ml/preprocessing/__init__.py](#)

[src/app/ml/preprocessing/data_prep.py](#)

"""Data preparation for ML training – fetches from MinIO Parquet files."""

```

import io
import logging
import os
from datetime import datetime
from typing import Optional, Tuple

```

```

import pandas as pd
import pyarrow.parquet as pq

logger = logging.getLogger(__name__)

FEATURE_COLS = ["temperature", "humidity", "light", "voltage"]
TIME_FEATURE_COLS = ["hour", "day_of_week"]
ALL_FEATURES = FEATURE_COLS + TIME_FEATURE_COLS

def _get_minio_client():
    from minio import Minio
    from src.app import config
    return Minio(
        config.MINIO_ENDPOINT,
        access_key=config.MINIO_ACCESS_KEY,
        secret_key=config.MINIO_SECRET_KEY,
        secure=config.MINIO_SECURE,
    )

def load_parquet_from_minio(
    bucket: Optional[str] = None,
    max_files: int = 0,
) -> pd.DataFrame:
    """Load Parquet files from MinIO and return a combined DataFrame.

    Think of MinIO as a filing cabinet where each drawer is a day's data.
    This function opens the drawers and reads all the files inside.

    Args:
        bucket: MinIO bucket name (defaults to config.MINIO_BUCKET)
        max_files: Maximum number of Parquet files to load (0 = all files)

    Returns:
        Combined DataFrame with all sensor readings
    """
    from src.app import config
    client = _get_minio_client()
    bucket = bucket or config.MINIO_BUCKET

    logger.info(f>Loading Parquet files from MinIO bucket={bucket}")
    objects = list(client.list_objects(bucket, recursive=True))
    parquet_objects = [o for o in objects if
o.object_name.endswith(".parquet")]

    if not parquet_objects:
        logger.warning("No Parquet files found in MinIO")
        return pd.DataFrame()

```

```

    if max_files > 0:
        logger.info(f"Found {len(parquet_objects)} Parquet files, loading up
to {max_files}")
        parquet_objects = parquet_objects[:max_files]
    else:
        logger.info(f"Found {len(parquet_objects)} Parquet files, loading
ALL")

frames = []
for obj in parquet_objects:
    try:
        response = client.get_object(bucket, obj.object_name)
        data = response.read()
        response.close()
        df = pd.read_parquet(io.BytesIO(data))
        frames.append(df)
    except Exception as e:
        logger.warning(f"Failed to load {obj.object_name}: {e}")

if not frames:
    logger.warning("All Parquet files failed to load")
    return pd.DataFrame()

combined = pd.concat(frames, ignore_index=True)
logger.info(f"Loaded {len(combined)} rows from {len(frames)} files")
return combined

def load_from_local_file(file_path: str, max_rows: int = 500_000) ->
pd.DataFrame:
    """Load sensor data from a local text file (historical_data.txt format).

    Columns: date, time, epoch, mote_id, temperature, humidity, light,
voltage, updated_timestamp
    """
    logger.info(f"Loading from local file: {file_path}")
    cols = ["date", "time", "epoch", "mote_id", "temperature", "humidity",
"light", "voltage", "updated_timestamp"]
    df = pd.read_csv(
        file_path,
        sep=r"\s+",
        names=cols,
        na_values=["N/A", "nan", ""],
        nrows=max_rows,
        on_bad_lines="skip",
    )
    # Combine date + time into timestamp
    try:
        df["timestamp"] = pd.to_datetime(df["date"] + " " + df["time"],

```

```

errors="coerce")
    except Exception:
        df["timestamp"] = pd.NaT
    return df

def prepare_features(df: pd.DataFrame) -> Tuple[pd.DataFrame, pd.Series]:
    """Clean and engineer features for ML training.

    Like a chef prepping ingredients – this removes bad data,
    adds time-based features, and returns a clean feature matrix.

    Returns:
        X: Feature DataFrame (ready for sklearn)
        mote_ids: Series of mote IDs aligned with X
    """
    df = df.copy()

    # Ensure sensor columns are numeric
    for col in FEATURE_COLS:
        if col in df.columns:
            df[col] = pd.to_numeric(df[col], errors="coerce")

    # Drop rows with all sensor values missing
    df = df.dropna(subset=FEATURE_COLS, how="all")

    # Parse timestamp columns from strings if needed
    # (consumer writes updated_timestamp as ISO strings in Parquet)
    for ts_col in ("updated_timestamp", "timestamp"):
        if ts_col in df.columns and not
pd.api.types.is_datetime64_any_dtype(df[ts_col]):
            df[ts_col] = pd.to_datetime(df[ts_col], utc=True, errors="coerce")

    # Add time features – prefer updated_timestamp (real-time mapped), fall
back to timestamp
    ts_series = None
    for ts_col in ("updated_timestamp", "timestamp"):
        if ts_col in df.columns and
pd.api.types.is_datetime64_any_dtype(df[ts_col]):
            ts_series = df[ts_col]
            break

    if ts_series is not None:
        df["hour"] = ts_series.dt.hour
        df["day_of_week"] = ts_series.dt.dayofweek
    else:
        df["hour"] = 12 # fallback
        df["day_of_week"] = 0

    # Keep only valid ranges (sensor physics bounds)

```

```

df = df[df["temperature"].between(-10, 80) | df["temperature"].isna()]
df = df[df["humidity"].between(0, 100) | df["humidity"].isna()]
df = df[df["light"].between(0, 3000) | df["light"].isna()]
df = df[df["voltage"].between(0, 5) | df["voltage"].isna()]

available_features = [c for c in ALL_FEATURES if c in df.columns]
X = df[available_features].copy()

# Fill remaining NaN with column medians
X = X.fillna(X.median(numeric_only=True))

mote_ids = df.get("mote_id", pd.Series(["unknown"] * len(df)))

logger.info(f"Prepared feature matrix: {X.shape},
features={available_features}")
return X, mote_ids

```

[src/app/scheduler/__init__.py](#)

SIEIS Scheduler package

[src/app/scheduler/jobs.py](#)

```

"""
SIEIS Scheduler Jobs
=====
Two daily jobs that automate the full pipeline:

    Job A   00:00 UTC   - Remap incremental_data.txt to today's dates
                        → Restart simulator so it emits fresh data
    Job B   02:00 UTC   - Retrain IsolationForest on yesterday's MinIO Parquet
files
                        → Hot-reload the FastAPI model (no restart needed)

```

Analogy

Think of this like a newspaper printing plant:

- Job A is the night-shift that updates the "today's date" stamp on the press
 - and restarts the press so it prints the right edition.
- Job B is the editor-in-chief who reviews all of yesterday's stories,
 - refines the spam filter (anomaly model), and pushes it live.

"""

```

import logging
import os
import sys
from datetime import date, timedelta
from pathlib import Path

```

```

import requests

logger = logging.getLogger(__name__)

# — Environment config


---


API_RELOAD_URL = os.getenv(
    "API_RELOAD_URL", "http://sieis-api:8000/api/v1/ml/model/reload"
)
SIMULATOR_CONTAINER = os.getenv("SIMULATOR_CONTAINER", "sieis-simulator")
INCR_DATA_PATH = Path(
    os.getenv("INCR_DATA_PATH", "/app/data/processed/incremental_data.txt")
)

# — Helpers


---



def _restart_simulator() -> bool:
    """Restart the simulator Docker container via the Docker Python SDK.

    Returns True on success, False on failure.

    The Docker socket (/var/run/docker.sock) must be mounted into this
    container (see docker-compose.yml).
    """
    try:
        import docker # docker>=6.1.0
        client = docker.from_env()
        container = client.containers.get(SIMULATOR_CONTAINER)
        container.restart(timeout=30)
        logger.info("Simulator container '%s' restarted successfully",
SIMULATOR_CONTAINER)
        return True
    except Exception:
        logger.exception(
            "Failed to restart simulator container '%s'. "
            "Check that /var/run/docker.sock is mounted and the container name
is correct.",
            SIMULATOR_CONTAINER,
        )
        return False

def _reload_api_model() -> bool:
    """POST to the FastAPI /ml/model/reload endpoint to hot-swap the model.

    Returns True on success, False on failure.
    No container restart required – the API picks up the latest .pkl from

```

```

disk.
    """
    try:
        resp = requests.post(API_RELOAD_URL, timeout=15)
        resp.raise_for_status()
        payload = resp.json()
        logger.info("API model reloaded: %s", payload)
        return True
    except requests.exceptions.ConnectionError:
        logger.error(
            "Could not reach API at %s. Is sieis-api running?", API_RELOAD_URL
        )
        return False
    except Exception:
        logger.exception("Unexpected error calling API reload endpoint")
        return False

```

— Job A: Daily Data Refresh

```

def data_exists_for_today() -> bool:
    """
    Check if there is data for today.
    You should implement this function to check your data source (e.g., MinIO,
    local file, DB).
    """
    from src.app.ml.preprocessing.data_prep import load_parquet_from_minio
    df = load_parquet_from_minio()
    today_str = date.today().isoformat()
    # Adjust the column name as per your data, e.g., 'date' or 'timestamp'
    return not df[df['date'] == today_str].empty

```

```

def remap_and_restart() -> None:
    """
    Modified Job A – runs at 00:00 UTC daily.
    Only restarts the simulator if today's data exists. No remapping.
    """
    logger.info("Restarting simulator container '%s'", SIMULATOR_CONTAINER)
    ok = _restart_simulator()
    status = "SUCCESS" if ok else "FAILED (check logs)"
    logger.info("Simulator restart: %s", status)

```

— Job B: Daily Model Retrain

```

def retrain_and_reload() -> None:
    """
    Job B – runs at 02:00 UTC daily.

    Step 1: Load yesterday's sensor data from MinIO Parquet archive.

```


Falls back to last 30 days of Parquet if yesterday alone is sparse.

Step 2: Prepare feature matrix: temperature, humidity, light, voltage, hour-of-day, day-of-week.

Step 3: Train a new IsolationForest model and save it as:
anomaly_detector_YYYYMMDD_HHMMSS_scheduled.pkl

Step 4: POST to FastAPI /ml/model/reload – the API hot-swaps its in-memory pipeline without any container restart.

Analogy: Like a bakery that checks what sold best yesterday, adjusts the recipe overnight, and opens in the morning with a fresher product – without closing the shop.

```
"""
logger.info("Loading Parquet data from MinIO (last 3 days)...")
try:
    from src.app.ml.preprocessing.data_prep import (
        load_parquet_from_minio,
        prepare_features,
    )
    from datetime import date, timedelta

    df = load_parquet_from_minio()

    today = date.today()
    recent_dates = [(today - timedelta(days=i)).isoformat() for i in
range(0, 3)]
    if 'date' in df.columns:
        df_recent = df[df['date'].isin(recent_dates)]
    else:
        logger.warning("DataFrame does not have a 'date' column. Using all
data.")
        df_recent = df

    if df_recent.empty:
        logger.warning(
            "No Parquet data found in MinIO for last 3 days. "
            "Retrain skipped – existing model remains active."
        )
        return

    logger.info("Loaded %d rows for retraining", len(df_recent))
except Exception:
    logger.exception("MinIO data load failed – retrain aborted")
    return

try:
    X, mote_ids = prepare_features(df_recent)
```

```

        if X.empty or len(X) < 100:
            logger.warning(
                "Feature matrix too small (%d rows) for reliable training. "
                "Need at least 100 rows. Retrain skipped.",
                len(X),
            )
            return
        logger.info("Feature matrix shape: %s", X.shape)
    except Exception:
        logger.exception("Feature preparation failed – retrain aborted")
        return

# -- Step 3: Train and save model -----
logger.info("[3/4] Training IsolationForest...")
try:
    from src.app.ml.detector import train_anomaly_detector, save_model #
noqa: PLC0415

    pipeline, metrics = train_anomaly_detector(X)
    logger.info(
        "Training complete: %d anomalies / %d samples (%.1f%%)",
        metrics["n_anomalies_detected"],
        metrics["n_samples"],
        metrics["anomaly_ratio"] * 100,
    )

    filename = save_model(pipeline, metrics, tag="scheduled")
    logger.info("Model saved: %s", filename)
except Exception:
    logger.exception("Model training/save failed – retrain aborted")
    return

# -- Step 4: Hot-reload API model -----
logger.info("[4/4] Reloading API model...")
ok = _reload_api_model()
status = "SUCCESS" if ok else "FAILED – API still uses previous model"
logger.info("JOB B complete – API reload: %s", status)
logger.info("=" * 60)

```

src/app/scheduler/main.py

```
"""
```

```
SIEIS Scheduler Entry Point
```

```
=====
```

```
Starts APScheduler with two daily cron jobs:
```

```

    Job A  00:00 UTC  remap_and_restart  – refresh incremental data + restart
simulator

```

```

    Job B  02:00 UTC  retrain_and_reload – retrain anomaly model + hot-reload

```

API

Environment Variables (all optional, with defaults):

| | | |
|--------------------|-------|---|
| SCHEDULER_TIMEZONE | UTC | |
| JOB_A_HOUR | 0 | (midnight UTC) |
| JOB_A_MINUTE | 0 | |
| JOB_B_HOUR | 2 | (2 AM UTC) |
| JOB_B_MINUTE | 0 | |
| RUN_JOBS_ON_START | false | set "true" to fire both jobs immediately at startup |

(useful for smoke-testing without waiting for midnight)

Usage:

```
docker-compose up scheduler          # normal operation
RUN_JOBS_ON_START=true docker-compose up scheduler  # immediate test run
"""
```

```
import logging
import os
import sys
from datetime import datetime
from pathlib import Path

# Ensure project root is on path when run inside Docker (/app)
_project_root = Path(__file__).parents[3] # src/app/scheduler -> project root
if str(_project_root) not in sys.path:
    sys.path.insert(0, str(_project_root))

from apscheduler.schedulers.blocking import BlockingScheduler
from apscheduler.triggers.cron import CronTrigger
from apscheduler.events import EVENT_JOB_ERROR, EVENT_JOB_EXECUTED

from src.app.scheduler.jobs import remap_and_restart, retrain_and_reload

logging.basicConfig(
    level=logging.INFO,
    format="%(asctime)s [%(levelname)s] %(name)s - %(message)s",
    datefmt="%Y-%m-%d %H:%M:%S",
)
logger = logging.getLogger(__name__)

# — Config from environment
TIMEZONE = os.getenv("SCHEDULER_TIMEZONE", "UTC")
JOB_A_HOUR = int(os.getenv("JOB_A_HOUR", "0"))
JOB_A_MINUTE = int(os.getenv("JOB_A_MINUTE", "0"))
JOB_B_HOUR = int(os.getenv("JOB_B_HOUR", "2"))
JOB_B_MINUTE = int(os.getenv("JOB_B_MINUTE", "0"))
RUN_JOBS_ON_START = os.getenv("RUN_JOBS_ON_START", "false").lower() == "true"
```

```
# — Event listener for job success / failure alerts
```

```
def _job_listener(event):
    """Log job outcomes clearly — makes docker logs easy to grep."""
    job = event.job_id
    if event.exception:
        logger.error("JOB FAILED [%s]: %s", job, event.exception)
    else:
        logger.info("JOB SUCCESS [%s]", job)
```

```
# — Main
```

```
def main():
    logger.info("=" * 60)
    logger.info("SIEIS Scheduler starting")
    logger.info(" Timezone : %s", TIMEZONE)
    logger.info(" Job A      : %02d:%02d  remap_and_restart", JOB_A_HOUR,
JOB_A_MINUTE)
    logger.info(" Job B      : %02d:%02d  retrain_and_reload", JOB_B_HOUR,
JOB_B_MINUTE)
    logger.info(" Immediate: %s", RUN_JOBS_ON_START)
    logger.info("=" * 60)

    scheduler = BlockingScheduler(timezone=TIMEZONE)
    scheduler.add_listener(_job_listener, EVENT_JOB_EXECUTED |
EVENT_JOB_ERROR)

    # — Job A: Daily data refresh —————
    scheduler.add_job(
        remap_and_restart,
        trigger=CronTrigger(hour=JOB_A_HOUR, minute=JOB_A_MINUTE,
timezone=TIMEZONE),
        id="job_a_data_refresh",
        name="Daily Data Refresh (remap + restart simulator)",
        max_instances=1,
        coalesce=True,          # if missed, run once (not multiple catch-ups)
        misfire_grace_time=300, # 5-min grace window before marking as missed
    )

    # — Job B: Daily model retrain —————
    scheduler.add_job(
        retrain_and_reload,
        trigger=CronTrigger(hour=JOB_B_HOUR, minute=JOB_B_MINUTE,
timezone=TIMEZONE),
        id="job_b_model_retrain",
```

```

        name="Daily Model Retrain (train + reload API)",
        max_instances=1,
        coalesce=True,
        misfire_grace_time=600, # 10-min grace (training can be slow)
    )

    # — Immediate run for smoke-testing —————
    if RUN_JOBS_ON_START:
        logger.info("RUN_JOBS_ON_START=true — triggering both jobs now for
testing")
        scheduler.add_job(
            remap_and_reload,
            id="job_a_immediate",
            name="Immediate Test: remap_and_reload",
        )
        scheduler.add_job(
            retrain_and_reload,
            id="job_b_immediate",
            name="Immediate Test: retrain_and_reload",
        )

    # — Start (blocks forever) —————
    logger.info("Scheduler running. Next runs:")
    try:
        scheduler.start()
    except (KeyboardInterrupt, SystemExit):
        logger.info("Scheduler stopped gracefully")

if __name__ == "__main__":
    main()

```

[src/app/simulator/__init__.py](#)

"""Simulator package for streaming Intel Lab sensor data."""

src/app/simulator/data_loader.py

```
"""Load and clean Intel Lab sensor data for simulation."""

import pandas as pd
import logging
from datetime import datetime
from typing import Dict, Optional, Tuple

logger = logging.getLogger(__name__)

def load_data_loader(data_path: str, mote_locs_path: str, filter_today_only:
bool = True) -> Tuple[Dict[int, pd.DataFrame], pd.DataFrame]:
    """
    Load Intel Lab sensor data and mote locations.

    Args:
        data_path: Path to data file (8-column legacy or 9-column with
updated_timestamp)
        mote_locs_path: Path to mote_locs.txt file
        filter_today_only: If True, only load records where updated_timestamp
is TODAY

    Returns:
        Tuple of (mote_data_dict, mote_locations_df)
        - mote_data_dict: {mote_id: DataFrame with sensor readings}
        - mote_locations_df: DataFrame with mote locations
    """

    logger.info(f"Loading sensor data from {data_path}")

    # Load sensor data
    # File columns: date time epoch moteid temperature humidity light voltage
[updated_timestamp]
    # 9th column (updated_timestamp) is optional for backward compatibility
    # Note: Only specify dtype for first 8 columns, let pandas infer the 9th
(if present)
    df = pd.read_csv(
        data_path,
        sep=r'\s+',
        header=None,
        names=['date', 'time', 'epoch', 'moteid', 'temperature', 'humidity',
'light', 'voltage', 'updated_timestamp'],
        na_values=['?']
    )

    # Ensure proper dtypes for numeric columns
    df['epoch'] = pd.to_numeric(df['epoch'], errors='coerce')
    df['moteid'] = pd.to_numeric(df['moteid'], errors='coerce')
```

```

df['temperature'] = pd.to_numeric(df['temperature'], errors='coerce')
df['humidity'] = pd.to_numeric(df['humidity'], errors='coerce')
df['light'] = pd.to_numeric(df['light'], errors='coerce')
df['voltage'] = pd.to_numeric(df['voltage'], errors='coerce')

raw_count = len(df)
logger.info(f"Loaded {raw_count:,} total records")

# — Row validation —————
# Shifted rows: when a source line has fewer fields than expected, pandas
# shifts all values left. The symptom is an ISO-8601 timestamp string
# landing in the `voltage` column and `updated_timestamp` becoming NaN.
# Detect these early by coercing numeric columns and flagging NaN moteids.

# Guard 1 – moteid must be numeric (already coerced above; NaN rows
dropped here)
non_numeric_mote = df['moteid'].isna().sum()
if non_numeric_mote:
    logger.warning(
        f"Row validation: dropped {non_numeric_mote:,} rows with non-
numeric moteid"
    )
df = df.dropna(subset=['moteid'])

# Guard 2 – moteid must be within the Intel Lab sensor range (1-100)
# Rows outside this range are artefacts of column-shift parse errors.
MOTE_ID_MIN, MOTE_ID_MAX = 1, 100
out_of_range = ~df['moteid'].between(MOTE_ID_MIN, MOTE_ID_MAX)
if out_of_range.any():
    logger.warning(
        f"Row validation: dropped {out_of_range.sum():,} rows with "
        f"moteid outside [{MOTE_ID_MIN}, {MOTE_ID_MAX}] "
        f"(saw values: {sorted(df.loc[out_of_range,
'moteid'].unique().tolist()[:10])}) "
    )
df = df[~out_of_range]

total_dropped = raw_count - len(df)
if total_dropped:
    logger.info(
        f"Row validation complete: {total_dropped:,} malformed rows
removed "
        f"({total_dropped / raw_count * 100:.2f}% of input), "
        f"{len(df):,} clean rows retained"
    )

# Combine date + time into a single timestamp column (original 2004
timestamps)
df['timestamp'] = pd.to_datetime(
    df['date'].astype(str) + ' ' + df['time'].astype(str), format='%Y-%m-

```

```

%d %H:%M:%S.%f', errors='coerce'
)

# Drop rows where timestamp could not be parsed
df = df.dropna(subset=['timestamp'])

# Parse updated_timestamp if present (pre-mapped to current year)
if 'updated_timestamp' in df.columns and
df['updated_timestamp'].notna().any():
    df['updated_timestamp'] = pd.to_datetime(df['updated_timestamp'],
errors='coerce')
    logger.info("Found updated_timestamp column in data file")

# Guard 3 – drop rows whose updated_timestamp failed to parse.
# In 9-column mode these are the remaining shifted rows: their voltage
# column holds a timestamp string, so `updated_timestamp` is NaN.
bad_ts = df['updated_timestamp'].isna()
if bad_ts.any():
    logger.warning(
        f"Row validation: dropped {bad_ts.sum():,} rows with "
        "unparseable updated_timestamp (likely shifted-column rows)"
    )
df = df[~bad_ts]

# FILTER: Only keep records for TODAY if requested
if filter_today_only:
    today = datetime.now().date()
    initial_count = len(df)
    df = df[df['updated_timestamp'].dt.date == today].copy()
    logger.info(f"Filtered to {len(df):,} records for today ({today})
from {initial_count:,} total records")

    if len(df) == 0:
        logger.warning(f"⚠️ No records found for today ({today})!")
        logger.warning("    This might mean:")
        logger.warning("    1. The transformation was run on a
different date")
        logger.warning("    2. All today's data was filtered as future
timestamps")
        logger.warning("    3. The realtime_data.txt file needs to be
regenerated")
    else:
        logger.info("No updated_timestamp column found - using legacy 8-column
format")

# Convert moteid to int (moteid is already validated and non-null above)
df['moteid'] = df['moteid'].astype(int)
# Drop rows with null temperature or humidity
df = df.dropna(subset=['temperature', 'humidity'])

```



```

# Fill null light with 0
df['light'] = df['light'].fillna(0.0)

# Forward fill voltage per mote to keep continuity in streams
df = df.sort_values(['moteid', 'timestamp'])
df['voltage'] = df.groupby('moteid')['voltage'].ffill()

# Fill any remaining null voltage with backward fill (for first readings)
df['voltage'] = df['voltage'].bfill()

# Split by mote ID for per-sensor simulation
mote_data_dict = {}
unique_motes = df['moteid'].unique()
logger.info(f"Processing {len(unique_motes)} unique motes")

for mote_id in unique_motes:
    mote_df = df[df['moteid'] == mote_id].copy()
    # Sort by timestamp
    mote_df = mote_df.sort_values('timestamp').reset_index(drop=True)
    mote_data_dict[mote_id] = mote_df
    logger.debug(f" Mote {mote_id}: {len(mote_df)} records")

# Load mote locations
mote_locs_df = pd.read_csv(
    mote_locs_path,
    sep=r'\s+',
    header=None,
    names=['moteid', 'x', 'y'],
    dtype={'moteid': int, 'x': float, 'y': float}
)

return mote_data_dict, mote_locs_df

def get_mote_data(data_path: str, mote_locs_path: str, mote_id: int) ->
pd.DataFrame:
    """
    Get sensor data for a specific mote.

    Args:
        data_path: Path to data.txt file
        mote_locs_path: Path to mote_locs.txt file
        mote_id: Mote ID to retrieve

    Returns:
        DataFrame with sensor readings for the mote
    """
    mote_data_dict, _ = load_data_loader(data_path, mote_locs_path)
    return mote_data_dict.get(mote_id, pd.DataFrame())

```

```

def get_mote_location(mote_locs_path: str, mote_id: int) ->
Optional[Tuple[float, float]]:
    """
    Get (x, y) location for a mote.

    Args:
        mote_locs_path: Path to mote_locs.txt file
        mote_id: Mote ID to retrieve

    Returns:
        Tuple of (x, y) coordinates
    """
    mote_locs_df = pd.read_csv(
        mote_locs_path,
        sep=r'\s+',
        header=None,
        names=['moteid', 'x', 'y'],
        dtype={'moteid': int, 'x': float, 'y': float}
    )

    loc = mote_locs_df[mote_locs_df['moteid'] == mote_id]
    if len(loc) > 0:
        return (loc.iloc[0]['x'], loc.iloc[0]['y'])
    return None

```

[src/app/simulator/emitter.py](#)

```

"""Emit per-mote readings to Kafka with time compression."""

```

```

import logging
import time
from typing import Optional
from src.app.config import YEAR_OFFSET

```

```

logger = logging.getLogger(__name__)

```

```

def emit_mote(mote_id: int, df, producer, speed_factor: float = 100.0,
stop_event: Optional[object] = None) -> None:

```

```

    """Emit rows for a single mote DataFrame to Kafka using the provided
    producer.

```

```

        - `df` is expected to contain columns: `timestamp`, `epoch`,
`temperature`,
        `humidity`, `light`, `voltage`, and optionally `updated_timestamp` (pre-
mapped).
        - `speed_factor` compresses real-time delays (delay_seconds = delta /
speed_factor).

```

```

- `stop_event` (optional) should be an object with `is_set()` method to
allow graceful shutdown.
"""
    prev_ts = None
    for _, row in df.iterrows():
        if stop_event is not None and getattr(stop_event, "is_set", lambda:
False)():
            logger.info("Stop event set for mote %s, exiting emitter",
mote_id)
            break

        ts = row["timestamp"]
        if prev_ts is not None:
            delta = (ts - prev_ts).total_seconds()
            if delta > 0:
                # Sleep scaled by speed_factor to simulate real-time pacing.
                delay = float(delta) / float(speed_factor)
                time.sleep(delay)

        # Use pre-calculated updated_timestamp if available, otherwise
        calculate at runtime
        import pandas as pd
        if 'updated_timestamp' in row.index and
pd.isna(row.get('updated_timestamp')):
            # Use pre-mapped timestamp from data file (preferred for accuracy)
            # Convert pandas Timestamp to datetime if needed
            updated_ts =
pd.to_datetime(row['updated_timestamp']).to_pydatetime()
        else:
            # Fallback: Calculate updated timestamp mapping 2004 data to
current year
            # Handle leap year edge case (Feb 29 in leap year -> Feb 28 in
non-leap year)
            try:
                updated_ts = ts.replace(year=ts.year + YEAR_OFFSET)
            except ValueError:
                # Feb 29 in leap year 2004 -> Feb 28 in non-leap year 2025
                updated_ts = ts.replace(year=ts.year + YEAR_OFFSET, day=28)

    msg = {
        "mote_id": int(mote_id),
        "timestamp": ts.isoformat(),
        "updated_timestamp": updated_ts.isoformat(),
        "original_timestamp": ts.strftime("%Y-%m-%dT%H:%M:%S.%f"),
        "temperature": float(row.get("temperature")) if
row.get("temperature") is not None else None,
        "humidity": float(row.get("humidity")) if row.get("humidity") is
not None else None,
        "light": float(row.get("light")) if row.get("light") is not None
else None,

```

```

        "voltage": float(row.get("voltage")) if row.get("voltage") is not
None else None,
        "epoch": int(row.get("epoch")) if row.get("epoch") is not None
else None,
    }

    try:
        producer.send(msg, key=mote_id)
    except Exception:
        logger.exception("Emitter failed sending message for mote %s",
mote_id)

    prev_ts = ts

```

[src/app/simulator/main.py](#)

```

"""Simulator entry point to run the orchestrator in the foreground."""

```

```

import logging
import time
from pathlib import Path
import sys

# Add project root to sys.path so imports work when run directly
_project_root = Path(__file__).parents[3]
if str(_project_root) not in sys.path:
    sys.path.insert(0, str(_project_root))

from src.app.simulator.orchestrator import Orchestrator

logger = logging.getLogger(__name__)

def main():
    logging.basicConfig(level=logging.INFO)
    orch = Orchestrator()
    try:
        orch.start()
        # Keep main thread alive while child threads run
        while True:
            time.sleep(1.0)
    except KeyboardInterrupt:
        logger.info("KeyboardInterrupt received - stopping")
        orch.stop()

if __name__ == "__main__":
    main()

```

src/app/simulator/orchestrator.py

```
"""Orchestrate per-mote emitter threads and producer lifecycle."""

import logging
import threading
from typing import Optional

from src.app.config import KAFKA_BROKER, KAFKA_TOPIC, SPEED_FACTOR, DATA_PATH,
MOTE_LOCS_PATH, FILTER_TODAY_ONLY
from src.app.simulator.data_loader import load_data_loader
from src.app.simulator.producer import Producer
from src.app.simulator.emitter import emit_mote

logger = logging.getLogger(__name__)

class Orchestrator:
    """Start per-mote emitters and manage graceful shutdown."""
    def __init__(self, broker: str = KAFKA_BROKER, topic: str = KAFKA_TOPIC,
max_motes: Optional[int] = None, filter_today_only: Optional[bool] = None):
        self.producer = Producer(broker=broker, topic=topic)
        # Use provided filter parameter, otherwise default to config value
        filter_param = filter_today_only if filter_today_only is not None else
FILTER_TODAY_ONLY
        self.mote_data, _ = load_data_loader(DATA_PATH, MOTE_LOCS_PATH,
filter_today_only=filter_param)
        self.max_motes = max_motes
        self.threads = []
        self.stop_event = threading.Event()

    def start(self):
        mote_ids = list(self.mote_data.keys())
        if self.max_motes:
            mote_ids = mote_ids[: self.max_motes]

        logger.info("Starting orchestrator for %d motes", len(mote_ids))
        for mote_id in mote_ids:
            # One thread per mote keeps each sensor stream independent.
            t = threading.Thread(target=self._run_mote, args=(mote_id,))
            t.daemon = True
            t.start()
            self.threads.append(t)

    def _run_mote(self, mote_id: int):
        df = self.mote_data[mote_id]
        try:
            emit_mote(mote_id, df, self.producer, SPEED_FACTOR,
stop_event=self.stop_event)
        except Exception:
```

```

        logger.exception("Orchestrator error for mote %s", mote_id)

    def stop(self):
        logger.info("Stopping orchestrator: signalling stop to emitters")
        self.stop_event.set()
        for t in self.threads:
            t.join(timeout=1.0)
        logger.info("Flushing producer and closing")
        self.producer.flush()
        self.producer.close()

```

src/app/simulator/producer.py

"""Kafka producer wrapper with JSON serialization and safe send behavior."""

```

import json
import logging
from typing import Any, Optional

```

```

from kafka import KafkaProducer

```

```

logger = logging.getLogger(__name__)

```

```

class Producer:

```

```

    """Simple Kafka producer wrapper using kafka-python.

```

```

    If the broker is unavailable, sends are logged and skipped.
    """

```

```

    def __init__(self, broker: str = "localhost:9092", topic: str =
"sensor_readings"):
        self.broker = broker
        self.topic = topic
        self._producer: Optional[KafkaProducer] = None
        try:
            self._producer = KafkaProducer(
                bootstrap_servers=[broker],
                value_serializer=lambda v: json.dumps(v).encode("utf-8"),
                key_serializer=lambda k: str(k).encode("utf-8") if k is not
None else None,
                retries=3,
                acks=1,
            )
            logger.info("Connected Kafka producer to %s", broker)
        except Exception as exc: # pragma: no cover - runtime environment
dependent
            logger.warning("Could not initialize KafkaProducer (%s); sending
will be no-op", exc)

```

```

        self._producer = None

    def send(self, value: Any, key: Any = None) -> None:
        """Send a message to the configured topic. Synchronous by default.

        This method will log and return if the producer was not initialized.
        """
        if not self._producer:
            logger.debug("Producer not initialized - skipping send: %s",
value)
            return

        try:
            fut = self._producer.send(self.topic, key=key, value=value)
            # Block until the broker acks the message.
            fut.get(timeout=10)
        except Exception:
            logger.exception("Failed to send message to Kafka: %s", value)

    def flush(self) -> None:
        if self._producer:
            try:
                self._producer.flush()
            except Exception:
                logger.exception("Kafka flush failed")

    def close(self) -> None:
        if self._producer:
            try:
                self._producer.close()
            except Exception:
                logger.exception("Kafka close failed")

```

tests/run_all_tests.py

```

"""Master test runner for all container tests and full pipeline."""

import sys
import subprocess
from pathlib import Path

def run_test(test_file):
    """Run a single test file and return result."""
    print(f"\n{'='*80}")
    print(f"Running: {test_file.name}")
    print(f"{'='*80}")

    result = subprocess.run(

```

```

        [sys.executable, str(test_file)],
        capture_output=False
    )

    return result.returncode == 0

def main():
    """Run all container tests and full pipeline test."""
    tests_dir = Path(__file__).parent

    tests = [
        tests_dir / "test_container_redpanda.py",
        tests_dir / "test_container_influxdb.py",
        tests_dir / "test_container_minio.py",
        tests_dir / "test_full_pipeline.py"
    ]

    print("\n" + "🚀" * 40)
    print("MASTER TEST SUITE - ALL CONTAINERS + FULL PIPELINE")
    print("🚀" * 40)
    print("\nThis will test:")
    print("  1. Redpanda (Kafka) container")
    print("  2. InfluxDB 2.7 container")
    print("  3. MinIO container")
    print("  4. Full end-to-end data pipeline")

    results = []
    for test_file in tests:
        if test_file.exists():
            passed = run_test(test_file)
            results.append((test_file.name, passed))
        else:
            print(f"\n⚠️ Test file not found: {test_file}")
            results.append((test_file.name, False))

    # Final summary
    print("\n" + "="*80)
    print("FINAL SUMMARY - ALL TESTS")
    print("="*80)

    for test_name, passed in results:
        status = "✅ PASS" if passed else "❌ FAIL"
        print(f"{status}: {test_name}")

    passed_count = sum(1 for _, p in results if p)
    total_count = len(results)

    print(f"\n🇪🇺 Overall: {passed_count}/{total_count} test suites passed")

```



```

if passed_count == total_count:
    print("\n🎉🎉🎉 ALL TESTS PASSED! 🎉🎉🎉")
    print("\n✅ All containers are operational")
    print("✅ Full data pipeline is working")
    print("✅ Dual-write pattern (InfluxDB + MinIO) verified")
    print("\nSystem is fully operational and ready for production use!")
    return 0
else:
    print(f"\n⚠️ {total_count - passed_count} test suite(s) failed")
    print("\n📋 Next steps:")
    print("  1. Review failed test output above")
    print("  2. Check Docker containers: docker ps")
    print("  3. Check logs: docker-compose logs -f")
    print("  4. Restart services if needed: docker-compose restart")
    return 1

if __name__ == "__main__":
    sys.exit(main())

```

tests/test_container_influxdb.py

"""Test InfluxDB 2.7 container functionality."""

```

import sys
from pathlib import Path
sys.path.insert(0, str(Path(__file__).parent.parent))

import time
from datetime import datetime, timedelta
from influxdb_client import InfluxDBClient, Point, WritePrecision
from influxdb_client.client.write_api import SYNCHRONOUS
from tests.test_config import INFLUX_URL, INFLUX_TOKEN, INFLUX_ORG,
INFLUX_BUCKET

def test_influxdb_container():
    """Test 1: Verify InfluxDB container is running and accessible."""
    print("\n" + "="*80)
    print("TEST 1: InfluxDB Container Health Check")
    print("="*80)

    try:
        client = InfluxDBClient(
            url=INFLUX_URL,
            token=INFLUX_TOKEN,
            org=INFLUX_ORG,

```

```

        timeout=10000
    )

    health = client.health()
    print(f"✅ Connected to InfluxDB at {INFLUX_URL}")
    print(f"✅ Health status: {health.status}")
    print(f"✅ Version: {health.version}")

    orgs_api = client.organizations_api()
    orgs = orgs_api.find_organizations()
    org_names = [org.name for org in orgs]
    print(f"✅ Organizations: {org_names}")

    if INFLUX_ORG in org_names:
        print(f"✅ Target organization found: {INFLUX_ORG}")

    client.close()
    return True
except Exception as e:
    print(f"❌ Failed to connect to InfluxDB: {e}")
    return False

def test_bucket_access():
    """Test 2: Verify bucket exists and is accessible."""
    print("\n" + "="*80)
    print("TEST 2: Bucket Access Verification")
    print("="*80)

    try:
        client = InfluxDBClient(
            url=INFLUX_URL,
            token=INFLUX_TOKEN,
            org=INFLUX_ORG
        )

        buckets_api = client.buckets_api()
        buckets = buckets_api.find_buckets().buckets
        bucket_names = [bucket.name for bucket in buckets]

        print(f"✅ Available buckets: {bucket_names}")

        if INFLUX_BUCKET in bucket_names:
            print(f"✅ Target bucket found: {INFLUX_BUCKET}")
            target_bucket = next(b for b in buckets if b.name ==
INFLUX_BUCKET)
            print(f"    - ID: {target_bucket.id}")
            print(f"    - Retention: {target_bucket.retention_rules}")

```

```

        client.close()
        return True
    else:
        print(f"❌ Target bucket not found: {INFLUX_BUCKET}")
        client.close()
        return False
except Exception as e:
    print(f"❌ Bucket access test failed: {e}")
    return False


def test_write_query():
    """Test 3: Verify write and query operations."""
    print("\n" + "="*80)
    print("TEST 3: Write and Query Operations")
    print("="*80)

    try:
        client = InfluxDBClient(
            url=INFLUX_URL,
            token=INFLUX_TOKEN,
            org=INFLUX_ORG
        )

        write_api = client.write_api(write_options=SYNCHRONOUS)

        test_points = []
        now = datetime.utcnow()

        for i in range(5):
            point = Point("test_measurement") \
                .tag("mote_id", "999") \
                .tag("test", "true") \
                .field("temperature", 20.0 + i * 0.5) \
                .field("humidity", 50.0 + i) \
                .time(now - timedelta(minutes=i), WritePrecision.NS)
            test_points.append(point)

        print(f"📝 Writing {len(test_points)} test points...")
        write_api.write(bucket=INFLUX_BUCKET, org=INFLUX_ORG,
            record=test_points)
        print(f"✅ Write successful")

        time.sleep(2)

        query_api = client.query_api()
        query = f'''
        from(bucket: "{INFLUX_BUCKET}")
        |> range(start: -1h)

```

```

        |> filter(fn: (r) => r._measurement == "test_measurement")
        |> filter(fn: (r) => r.mote_id == "999")
'''

print("🔍 Executing Flux query...")
tables = query_api.query(query, org=INFLUX_ORG)

record_count = 0
for table in tables:
    for record in table.records:
        record_count += 1
        if record_count <= 3:
            print(f"✅ Record: {record.get_field()} = {record.get_value()}")

print(f"✅ Query returned {record_count} records")
client.close()
return record_count > 0
except Exception as e:
    print(f"❌ Write/Query test failed: {e}")
    return False

def test_sensor_reading_check():
    """Test 4: Check if sensor_reading measurement has data."""
    print("\n" + "="*80)
    print("TEST 4: Sensor Reading Measurement Check")
    print("="*80)

    try:
        client = InfluxDBClient(
            url=INFLUX_URL,
            token=INFLUX_TOKEN,
            org=INFLUX_ORG
        )

        query_api = client.query_api()
        query = f'''
from(bucket: "{INFLUX_BUCKET}")
  |> range(start: -24h)
  |> filter(fn: (r) => r._measurement == "sensor_reading")
  |> limit(n: 5)
'''

        print("🔍 Checking for sensor_reading data...")
        tables = query_api.query(query, org=INFLUX_ORG)

        record_count = 0
        mote_ids = set()

```

```

        for table in tables:
            for record in table.records:
                record_count += 1
                mote_id = record.values.get("mote_id")
                if mote_id:
                    mote_ids.add(mote_id)
                if record_count <= 3:
                    print(f" ✅ Sample: mote={mote_id},
{record.get_field()}={record.get_value():.2f}")

            if record_count > 0:
                print(f" ✅ Found {record_count} sensor records from
{len(mote_ids)} motes")
            else:
                print(f" ⚠️ No sensor_reading data found yet (simulator may not
have run)")

        client.close()
        return True
    except Exception as e:
        print(f" ❌ Sensor data check failed: {e}")
        return False

def main():
    """Run all InfluxDB container tests."""
    print("\n" + "📦" * 40)
    print("INFLUXDB 2.7 CONTAINER TEST SUITE")
    print("📦" * 40)
    print(f"\nTarget: {INFLUX_URL}")
    print(f"Organization: {INFLUX_ORG}")
    print(f"Bucket: {INFLUX_BUCKET}")

    results = []

    results.append(("Container Health", test_influxdb_container()))
    results.append(("Bucket Access", test_bucket_access()))
    results.append(("Write/Query", test_write_query()))
    results.append(("Sensor Data Check", test_sensor_reading_check()))

    print("\n" + "="*80)
    print("TEST SUMMARY")
    print("="*80)

    passed = sum(1 for _, result in results if result)
    total = len(results)

    for test_name, result in results:

```

```

        status = "✅ PASS" if result else "❌ FAIL"
        print(f"{status}: {test_name}")

    print(f"\nTotal: {passed}/{total} tests passed")

    if passed == total:
        print("\n🎉 All InfluxDB tests passed!")
        return 0
    else:
        print(f"\n⚠️ {total - passed} test(s) failed")
        return 1

if __name__ == "__main__":
    sys.exit(main())

```

tests/test_container_minio.py

```

"""Test MinIO container functionality."""

import sys
from pathlib import Path
sys.path.insert(0, str(Path(__file__).parent.parent))

import io
import time
from datetime import datetime
from minio import Minio
from tests.test_config import (
    MINIO_ENDPOINT,
    MINIO_ACCESS_KEY,
    MINIO_SECRET_KEY,
    MINIO_BUCKET,
    MINIO_SECURE
)

def test_minio_container():
    """Test 1: Verify MinIO container is running and accessible."""
    print("\n" + "="*80)
    print("TEST 1: MinIO Container Health Check")
    print("="*80)

    try:
        client = Minio(
            endpoint=MINIO_ENDPOINT,
            access_key=MINIO_ACCESS_KEY,
            secret_key=MINIO_SECRET_KEY,
            secure=MINIO_SECURE

```

```

    )

    buckets = client.list_buckets()
    bucket_names = [bucket.name for bucket in buckets]

    print(f"✅ Connected to MinIO at {MINIO_ENDPOINT}")
    print(f"✅ Available buckets: {bucket_names}")

    return True
except Exception as e:
    print(f"❌ Failed to connect to MinIO: {e}")
    return False

def test_bucket_exists():
    """Test 2: Verify target bucket exists."""
    print("\n" + "="*80)
    print("TEST 2: Target Bucket Verification")
    print("="*80)

    try:
        client = Minio(
            endpoint=MINIO_ENDPOINT,
            access_key=MINIO_ACCESS_KEY,
            secret_key=MINIO_SECRET_KEY,
            secure=MINIO_SECURE
        )

        if client.bucket_exists(MINIO_BUCKET):
            print(f"✅ Target bucket exists: {MINIO_BUCKET}")

            try:
                objects = []
                for obj in client.list_objects(MINIO_BUCKET, recursive=True):
                    objects.append(obj)
                    if len(objects) >= 10:
                        break
                print(f"✅ Bucket contains {len(objects)} objects (showing max
10)")

                for obj in objects[:3]:
                    print(f"  - {obj.object_name} ({obj.size} bytes)")
            except Exception as e:
                print(f"⚠️ Could not list objects: {e}")

            return True
        else:
            print(f"⚠️ Target bucket does not exist: {MINIO_BUCKET}")

```

```

    print(" Attempting to create...")

    try:
        client.make_bucket(MINIO_BUCKET)
        print(f"✅ Created bucket: {MINIO_BUCKET}")
        return True
    except Exception as e:
        print(f"❌ Failed to create bucket: {e}")
        return False
except Exception as e:
    print(f"❌ Bucket check failed: {e}")
    return False

def test_upload_download():
    """Test 3: Verify upload and download operations."""
    print("\n" + "="*80)
    print("TEST 3: Upload/Download Operations")
    print("="*80)

    try:
        client = Minio(
            endpoint=MINIO_ENDPOINT,
            access_key=MINIO_ACCESS_KEY,
            secret_key=MINIO_SECRET_KEY,
            secure=MINIO_SECURE
        )

        test_object = f"test/upload_test_{int(time.time())}.txt"
        test_data = f"Test upload at {datetime.now().isoformat()}\n" * 10
        test_bytes = test_data.encode('utf-8')

        print(f"📁 Uploading test object: {test_object}")
        client.put_object(
            bucket_name=MINIO_BUCKET,
            object_name=test_object,
            data=io.BytesIO(test_bytes),
            length=len(test_bytes),
            content_type='text/plain'
        )
        print(f"✅ Upload successful")

        print(f"📁 Downloading test object...")
        response = client.get_object(MINIO_BUCKET, test_object)
        downloaded_data = response.read()
        response.close()
        response.release_conn()

        if downloaded_data == test_bytes:

```



```

        print("✅ Download successful, data matches")
    else:
        print("⚠️ Downloaded data does not match uploaded data")
        return False

    client.remove_object(MINIO_BUCKET, test_object)
    print("✅ Test object cleaned up")

    return True
except Exception as e:
    print(f"❌ Upload/Download test failed: {e}")
    return False

def test_parquet_structure():
    """Test 4: Verify Parquet file structure exists."""
    print("\n" + "="*80)
    print("TEST 4: Parquet File Structure Check")
    print("="*80)

    try:
        client = Minio(
            endpoint=MINIO_ENDPOINT,
            access_key=MINIO_ACCESS_KEY,
            secret_key=MINIO_SECRET_KEY,
            secure=MINIO_SECURE
        )

        print("🔍 Looking for year-partitioned Parquet files...")

        year_objects = list(client.list_objects(MINIO_BUCKET, prefix="year=",
            recursive=True))

        if year_objects:
            print(f"✅ Found {len(year_objects)} partitioned files")

            parquet_files = [obj for obj in year_objects if
                obj.object_name.endswith('.parquet')]
            print(f"✅ Parquet files: {len(parquet_files)}")

            if parquet_files:
                print("\n📁 Sample file structure:")
                for obj in parquet_files[:5]:
                    print(f"    - {obj.object_name}")
                    print(f"        Size: {obj.size/1024:.2f} KB, Modified: {obj.last_modified}")

            return True

    except Exception as e:
        print(f"❌ Parquet file structure test failed: {e}")
        return False

```

```

        else:
            print("⚠ No year-partitioned files found yet")
            print("    (Consumer may not have written data yet)")
            return True
    except Exception as e:
        print(f"❌ Parquet structure check failed: {e}")
        return False

def main():
    """Run all MinIO container tests."""
    print("\n" + "═ " * 40)
    print("MINIO CONTAINER TEST SUITE")
    print("═ " * 40)
    print(f"\nTarget: {MINIO_ENDPOINT}")
    print(f"Bucket: {MINIO_BUCKET}")

    results = []

    results.append(("Container Health", test_minio_container()))
    results.append(("Bucket Verification", test_bucket_exists()))
    results.append(("Upload/Download", test_upload_download()))
    results.append(("Parquet Structure", test_parquet_structure()))

    print("\n" + "═"*80)
    print("TEST SUMMARY")
    print("═"*80)

    passed = sum(1 for _, result in results if result)
    total = len(results)

    for test_name, result in results:
        status = "✅ PASS" if result else "❌ FAIL"
        print(f"{status}: {test_name}")

    print(f"\nTotal: {passed}/{total} tests passed")

    if passed == total:
        print("\n🎉 All MinIO tests passed!")
        return 0
    else:
        print(f"\n⚠ {total - passed} test(s) failed")
        return 1

if __name__ == "__main__":
    sys.exit(main())

```

tests/test_container_redpanda.py

```
"""Test Redpanda (Kafka) container functionality."""
```

```
import sys
from pathlib import Path
sys.path.insert(0, str(Path(__file__).parent.parent))

import json
import time
from kafka import KafkaProducer, KafkaConsumer
from kafka.admin import KafkaAdminClient, NewTopic
from tests.test_config import KAFKA_BROKER, KAFKA_TOPIC
```

```
def test_redpanda_container():
    """Test 1: Verify Redpanda container is running and accessible."""
    print("\n" + "="*80)
    print("TEST 1: Redpanda Container Health Check")
    print("="*80)

    try:
        admin = KafkaAdminClient(
            bootstrap_servers=[KAFKA_BROKER],
            client_id='test_admin',
            request_timeout_ms=5000
        )

        cluster_metadata = admin.list_topics()
        print(f"✅ Connected to Redpanda at {KAFKA_BROKER}")
        print(f"✅ Available topics: {list(cluster_metadata)}")

        admin.close()
        return True
    except Exception as e:
        print(f"❌ Failed to connect to Redpanda: {e}")
        return False
```

```
def test_topic_creation():
    """Test 2: Verify topic creation and management."""
    print("\n" + "="*80)
    print("TEST 2: Topic Creation and Management")
    print("="*80)

    try:
        admin = KafkaAdminClient(
            bootstrap_servers=[KAFKA_BROKER],
            request_timeout_ms=5000
```

```

    )

    test_topic = "test_sensor_readings"
    existing_topics = admin.list_topics()

    if test_topic not in existing_topics:
        topic = NewTopic(
            name=test_topic,
            num_partitions=3,
            replication_factor=1
        )
        admin.create_topics([topic])
        print(f"✅ Created topic: {test_topic}")
    else:
        print(f"✅ Topic already exists: {test_topic}")

    if KAFKA_TOPIC in existing_topics:
        print(f"✅ Main topic exists: {KAFKA_TOPIC}")
    else:
        print(f"⚠️ Main topic not found: {KAFKA_TOPIC} (will be auto-
created)")

    admin.close()
    return True
except Exception as e:
    print(f"❌ Topic management failed: {e}")
    return False

def test_producer_consumer():
    """Test 3: Verify message production and consumption."""
    print("\n" + "="*80)
    print("TEST 3: Message Production and Consumption")
    print("="*80)

    test_topic = "test_sensor_readings"
    test_messages = [
        {"mote_id": 999, "temperature": 20.5, "timestamp": "2025-02-
16T10:00:00"},
        {"mote_id": 999, "temperature": 21.0, "timestamp": "2025-02-
16T10:01:00"},
    ]

    try:
        producer = KafkaProducer(
            bootstrap_servers=[KAFKA_BROKER],
            value_serializer=lambda v: json.dumps(v).encode('utf-8'),
            key_serializer=lambda k: str(k).encode('utf-8') if k is not None
        else None,

```

```

        request_timeout_ms=5000
    )

    print(f"🚀 Sending {len(test_messages)} test messages...")
    for msg in test_messages:
        future = producer.send(test_topic, value=msg)
        future.get(timeout=5)
        print(f"✅ Sent: {msg}")

    producer.flush()
    producer.close()

    print(f"\n📡 Consuming messages from {test_topic}...")
    consumer = KafkaConsumer(
        test_topic,
        bootstrap_servers=[KAFKA_BROKER],
        auto_offset_reset='earliest',
        consumer_timeout_ms=5000,
        value_deserializer=lambda m: json.loads(m.decode('utf-8'))
    )

    received_messages = []
    for message in consumer:
        received_messages.append(message.value)
        print(f"✅ Received: {message.value}")
        if len(received_messages) >= len(test_messages):
            break

    consumer.close()

    if len(received_messages) >= len(test_messages):
        print(f"\n✅ Successfully sent and received {len(test_messages)}
messages")
        return True
    else:
        print(f"\n⚠️ Sent {len(test_messages)}, received
{len(received_messages)}")
        return False
    except Exception as e:
        print(f"❌ Producer/Consumer test failed: {e}")
        return False

def main():
    """Run all Redpanda container tests."""
    print("\n" + "🔥" * 40)
    print("REDPANDA (KAFKA) CONTAINER TEST SUITE")
    print("🔥" * 40)

```

```

print(f"\nTarget: {KAFKA_BROKER}")
print(f"Topic: {KAFKA_TOPIC}")

results = []

results.append(("Container Health", test_redpanda_container()))
results.append(("Topic Management", test_topic_creation()))
results.append(("Producer/Consumer", test_producer_consumer()))

print("\n" + "="*80)
print("TEST SUMMARY")
print("="*80)

passed = sum(1 for _, result in results if result)
total = len(results)

for test_name, result in results:
    status = "✅ PASS" if result else "❌ FAIL"
    print(f"{status}: {test_name}")

print(f"\nTotal: {passed}/{total} tests passed")

if passed == total:
    print("\n🎉 All Redpanda tests passed!")
    return 0
else:
    print(f"\n⚠️ {total - passed} test(s) failed")
    return 1

if __name__ == "__main__":
    sys.exit(main())

```

tests/test_data_loader.py

```

"""Tests for Intel Lab data loading and cleaning."""

import sys
from pathlib import Path
sys.path.insert(0, str(Path(__file__).parent.parent))

import pytest
from src.app.simulator.data_loader import load_data_loader

class TestDataLoader:
    """Test cases for data loader module."""

    @pytest.fixture(scope="class")

```

```

def data_paths(self):
    """Fixture providing paths to test data files."""
    return {
        'data': 'data/raw/data.txt',
        'mote_locs': 'data/raw/mote_locs.txt'
    }

def test_load_without_error(self, data_paths):
    """Test 1: File loads without error."""
    try:
        mote_data, mote_locs = load_data_loader(data_paths['data'],
data_paths['mote_locs'])
        assert True, "Data loaded successfully"
    except Exception as e:
        pytest.fail(f"Data loader failed with error: {e}")

def test_returns_mote_groups(self, data_paths):
    """Test 2: Returns mote groups (at least 1, typically around 54)."""
    mote_data, mote_locs = load_data_loader(data_paths['data'],
data_paths['mote_locs'])

    # Check that we have mote groups
    assert isinstance(mote_data, dict), "Mote data should be a dictionary"
    assert len(mote_data) > 0, "Should have at least one mote group"

    # Typically Intel Lab has 54 motes, but some may be missing
    # Accept 40-60 motes as valid
    assert 30 <= len(mote_data) <= 100, f"Expected 30-100 motes, got
{len(mote_data)}"

def test_each_group_sorted_by_timestamp(self, data_paths):
    """Test 3: Each group is sorted by timestamp."""
    mote_data, _ = load_data_loader(data_paths['data'],
data_paths['mote_locs'])

    for mote_id, df in mote_data.items():
        # Check that timestamps are in order
        timestamps = df['timestamp'].values
        assert (timestamps[:-1] <= timestamps[1:]).all(), \
            f"Mote {mote_id} is not sorted by timestamp"

def test_no_null_temperature(self, data_paths):
    """Test 4: No null values in temperature column."""
    mote_data, _ = load_data_loader(data_paths['data'],
data_paths['mote_locs'])

    for mote_id, df in mote_data.items():
        assert df['temperature'].isnull().sum() == 0, \
            f"Mote {mote_id} has null temperature values"

```

```

def test_no_null_humidity(self, data_paths):
    """Test 5: No null values in humidity column."""
    mote_data, _ = load_data_loader(data_paths['data'],
data_paths['mote_locs'])

    for mote_id, df in mote_data.items():
        assert df['humidity'].isnull().sum() == 0, \
            f"Mote {mote_id} has null humidity values"

```

tests/test_e2e_pipeline.py

```
"""
```

Full SIEIS Pipeline Test Suite

Tests the complete data flow:

1. Docker containers running
2. Data in InfluxDB (incremental only)
3. Data in MinIO (historical + incremental)
4. API endpoints working
5. ML model inference
6. Dashboard data loading
7. End-to-end validation

Usage:

```

python scripts/test_full_pipeline.py
python scripts/test_full_pipeline.py --verbose
python scripts/test_full_pipeline.py --skip-api # Skip API tests
"""

```

```

import sys
from pathlib import Path
sys.path.insert(0, str(Path(__file__).parent.parent))

import argparse
import logging
import subprocess
from datetime import datetime, timedelta
from time import sleep

logging.basicConfig(
    level=logging.INFO,
    format='%(asctime)s - %(levelname)s - %(message)s'
)
logger = logging.getLogger(__name__)

```

```

class FullPipelineTest:
    """Complete pipeline validation test."""

```



```

def __init__(self, verbose=False, skip_api=False):
    self.verbose = verbose
    self.skip_api = skip_api
    self.test_results = {}
    self.failed_tests = []
    self.passed_tests = []

# ===== PHASE 1: DOCKER TESTS =====

def test_1_docker_containers_running(self):
    """TEST 1: All Docker containers are running."""
    logger.info("\n" + "="*80)
    logger.info("TEST 1: Docker Containers Running")
    logger.info("="*80)

    try:
        result = subprocess.run(
            ["docker", "ps", "--filter", "name=sieis"],
            capture_output=True,
            text=True,
            timeout=10
        )

        required_containers = [
            "sieis-simulator",
            "sieis-consumer",
            "sieis-influxdb3",
            "sieis-minio",
            "sieis-redpanda"
        ]

        missing = []
        for container in required_containers:
            if container not in result.stdout:
                missing.append(container)

        if missing:
            logger.error(f"❌ FAIL: Missing containers: {missing}")
            self.failed_tests.append("Docker containers")
            return False

        logger.info(f"✅ PASS: All required containers running")
        for container in required_containers:
            logger.info(f"    ✅ {container}")

        self.passed_tests.append("Docker containers")
        return True

    except Exception as e:

```

```

        logger.error(f"❌ FAIL: {e}")
        self.failed_tests.append("Docker containers")
        return False

def test_2_influxdb_connectivity(self):
    """TEST 2: InfluxDB is accessible."""
    logger.info("\n" + "="*80)
    logger.info("TEST 2: InfluxDB Connectivity")
    logger.info("="*80)

    try:
        result = subprocess.run(
            ["curl", "-s", "http://localhost:8086/health"],
            capture_output=True,
            text=True,
            timeout=10
        )

        if '"status": "pass"' in result.stdout or "healthy" in result.stdout:
            logger.info("✅ PASS: InfluxDB is healthy")
            self.passed_tests.append("InfluxDB connectivity")
            return True
        else:
            logger.error(f"❌ FAIL: InfluxDB not healthy. Response: {result.stdout}")
            self.failed_tests.append("InfluxDB connectivity")
            return False

    except Exception as e:
        logger.error(f"❌ FAIL: {e}")
        self.failed_tests.append("InfluxDB connectivity")
        return False

def test_3_minio_connectivity(self):
    """TEST 3: MinIO is accessible."""
    logger.info("\n" + "="*80)
    logger.info("TEST 3: MinIO Connectivity")
    logger.info("="*80)

    try:
        from minio import Minio
        from src.app.config import MINIO_ENDPOINT, MINIO_ACCESS_KEY, MINIO_SECRET_KEY

        client = Minio(
            MINIO_ENDPOINT,
            access_key=MINIO_ACCESS_KEY,
            secret_key=MINIO_SECRET_KEY

```

```

    )

    # Try to list buckets
    buckets = client.list_buckets()
    bucket_names = [b.name for b in buckets.buckets]

    if "sieis-archive" in bucket_names:
        logger.info("✅ PASS: MinIO is accessible with 'sieis-archive'
bucket")

        self.passed_tests.append("MinIO connectivity")
        return True
    else:
        logger.error(f"❌ FAIL: 'sieis-archive' bucket not found")
        self.failed_tests.append("MinIO connectivity")
        return False

except Exception as e:
    logger.error(f"❌ FAIL: {e}")
    self.failed_tests.append("MinIO connectivity")
    return False

# ===== PHASE 2: DATA QUALITY TESTS =====

def test_4_influxdb_has_incremental_data(self):
    """TEST 4: InfluxDB has incremental data (last 24h)."""
    logger.info("\n" + "="*80)
    logger.info("TEST 4: InfluxDB Has Incremental Data")
    logger.info("="*80)

    try:
        from influxdb_client import InfluxDBClient
        from src.app.config import INFLUX_URL, INFLUX_TOKEN, INFLUX_ORG,
INFLUX_BUCKET

        client = InfluxDBClient(
            url=INFLUX_URL,
            token=INFLUX_TOKEN,
            org=INFLUX_ORG
        )
        query_api = client.query_api()

        query = f'''from(bucket: "{INFLUX_BUCKET}")
            |> range(start: -24h)
            |> filter(fn: (r) => r._measurement == "sensor_reading")
            |> count()
        '''

        result = query_api.query(query)

```

```

count = 0
for table in result:
    for record in table.records:
        count += record.get_value()

client.close()

if count > 0:
    logger.info(f"✅ PASS: Found {count:,} records in last 24h")
    self.passed_tests.append("InfluxDB incremental data")
    return True
else:
    logger.error(f"❌ FAIL: No data found in InfluxDB (last 24h)")
    self.failed_tests.append("InfluxDB incremental data")
    return False

except Exception as e:
    logger.error(f"❌ FAIL: {e}")
    self.failed_tests.append("InfluxDB incremental data")
    return False

def test_5_mote_count_correct(self):
    """TEST 5: Mote count is 42-44."""
    logger.info("\n" + "="*80)
    logger.info("TEST 5: Correct Mote Count (42-44)")
    logger.info("="*80)

    try:
        from influxdb_client import InfluxDBClient
        from src.app.config import INFLUX_URL, INFLUX_TOKEN, INFLUX_ORG,
INFLUX_BUCKET

        client = InfluxDBClient(
            url=INFLUX_URL,
            token=INFLUX_TOKEN,
            org=INFLUX_ORG
        )
        query_api = client.query_api()

        query = f'''from(bucket: "{INFLUX_BUCKET}")
            |> range(start: -24h)
            |> filter(fn: (r) => r._measurement == "sensor_reading")
            |> group(columns: ["mote_id"])
            |> distinct(column: "mote_id")
            |> count()
        '''

        result = query_api.query(query)

```

```

        mote_count = 0
        for table in result:
            for record in table.records:
                mote_count += 1

        client.close()

        if 42 <= mote_count <= 44:
            logger.info(f"✅ PASS: Found {mote_count} motes (expected 42-
44)")

            self.passed_tests.append("Mote count")
            return True
        else:
            logger.error(f"❌ FAIL: Found {mote_count} motes (expected 42-
44)")

            self.failed_tests.append("Mote count")
            return False

    except Exception as e:
        logger.error(f"❌ FAIL: {e}")
        self.failed_tests.append("Mote count")
        return False

def test_6_sensor_values_valid(self):
    """TEST 6: Sensor values are within valid ranges."""
    logger.info("\n" + "="*80)
    logger.info("TEST 6: Sensor Values Within Valid Ranges")
    logger.info("="*80)

    try:
        from influxdb_client import InfluxDBClient
        from src.app.config import INFLUX_URL, INFLUX_TOKEN, INFLUX_ORG,
INFLUX_BUCKET

        client = InfluxDBClient(
            url=INFLUX_URL,
            token=INFLUX_TOKEN,
            org=INFLUX_ORG
        )
        query_api = client.query_api()

        valid_ranges = {
            'temperature': (-10, 50),
            'humidity': (0, 100),
            'light': (0, 5000),
            'voltage': (2.0, 3.5),
        }

        all_valid = True

```

```

for field, (min_val, max_val) in valid_ranges.items():
    query = f'''from(bucket: "{INFLUX_BUCKET}")
        |> range(start: -24h)
        |> filter(fn: (r) => r._measurement == "sensor_reading")
        |> filter(fn: (r) => r._field == "{field}")
    '''

    result = query_api.query(query)

    out_of_range = 0
    for table in result:
        for record in table.records:
            value = record.get_value()
            if value is not None and not (min_val <= value <=
max_val):
                out_of_range += 1

    if out_of_range > 0:
        logger.warning(f"⚠️ {field}: {out_of_range} out-of-range
values")
        all_valid = False
    else:
        logger.info(f"✅ {field}: All values in range ({min_val}-
{max_val}))")

    client.close()

    if all_valid:
        logger.info("✅ PASS: All sensor values valid")
        self.passed_tests.append("Sensor values")
        return True
    else:
        logger.warning("⚠️ PASS WITH WARNINGS: Some out-of-range
values")
        self.passed_tests.append("Sensor values")
        return True

except Exception as e:
    logger.error(f"❌ FAIL: {e}")
    self.failed_tests.append("Sensor values")
    return False

def test_7_minio_has_data(self):
    """TEST 7: MinIO has Parquet data files."""
    logger.info("\n" + "="*80)
    logger.info("TEST 7: MinIO Has Data Files")
    logger.info("="*80)

    try:

```

```

        from minio import Minio
        from src.app.config import MINIO_ENDPOINT, MINIO_ACCESS_KEY,
MINIO_SECRET_KEY

        client = Minio(
            MINIO_ENDPOINT,
            access_key=MINIO_ACCESS_KEY,
            secret_key=MINIO_SECRET_KEY
        )

        # Count parquet files
        objects = client.list_objects("sieis-archive", recursive=True)
        parquet_count = 0
        file_list = []

        for obj in objects:
            if obj.object_name.endswith('.parquet'):
                parquet_count += 1
                file_list.append(obj.object_name)
                if len(file_list) <= 5:
                    logger.info(f"    Sample: {obj.object_name}")

        if parquet_count > 0:
            logger.info(f"✅ PASS: Found {parquet_count} Parquet files in
MinIO")

            self.passed_tests.append("MinIO data files")
            return True
        else:
            logger.error(f"❌ FAIL: No Parquet files found in MinIO")
            self.failed_tests.append("MinIO data files")
            return False

    except Exception as e:
        logger.error(f"❌ FAIL: {e}")
        self.failed_tests.append("MinIO data files")
        return False

# ===== PHASE 3: RECENT DATA TESTS =====

def test_8_recent_data_within_1_hour(self):
    """TEST 8: Most recent data is within 1 hour."""
    logger.info("\n" + "="*80)
    logger.info("TEST 8: Recent Data Within 1 Hour")
    logger.info("="*80)

    try:
        from influxdb_client import InfluxDBClient
        from src.app.config import INFLUX_URL, INFLUX_TOKEN, INFLUX_ORG,
INFLUX_BUCKET

```

```

client = InfluxDBClient(
    url=INFLUX_URL,
    token=INFLUX_TOKEN,
    org=INFLUX_ORG
)
query_api = client.query_api()

query = f'''from(bucket: "{INFLUX_BUCKET}")
  |> range(start: -24h)
  |> filter(fn: (r) => r._measurement == "sensor_reading")
  |> sort(columns: ["_time"], desc: true)
  |> limit(n: 1)
'''

result = query_api.query(query)

latest = None
for table in result:
    for record in table.records:
        latest = record.get_time()
        break

client.close()

if latest:
    now = datetime.now(latest.tzinfo)
    diff_minutes = (now - latest).total_seconds() / 60

    if diff_minutes < 60:
        logger.info(f"✅ PASS: Latest data is {diff_minutes:.1f}
minutes old")
        self.passed_tests.append("Recent data freshness")
        return True
    else:
        logger.error(f"❌ FAIL: Latest data is {diff_minutes:.1f}
minutes old (expected < 60)")
        self.failed_tests.append("Recent data freshness")
        return False
else:
    logger.error(f"❌ FAIL: Could not determine latest timestamp")
    self.failed_tests.append("Recent data freshness")
    return False

except Exception as e:
    logger.error(f"❌ FAIL: {e}")
    self.failed_tests.append("Recent data freshness")
    return False

```



```

# ===== PHASE 4: API TESTS =====

def test_9_api_health_endpoint(self):
    """TEST 9: API /health endpoint responds."""
    if self.skip_api:
        logger.info("\n🔊 Skipping API tests (--skip-api flag)")
        return None

    logger.info("\n" + "="*80)
    logger.info("TEST 9: API Health Endpoint")
    logger.info("="*80)

    try:
        result = subprocess.run(
            ["curl", "-s", "http://localhost:8000/health"],
            capture_output=True,
            text=True,
            timeout=10
        )

        if '"status": "healthy"' in result.stdout or "healthy" in result.stdout:
            logger.info("✅ PASS: API health endpoint responds")
            self.passed_tests.append("API health")
            return True
        else:
            logger.error(f"❌ FAIL: API not responding. Response: {result.stdout}")
            self.failed_tests.append("API health")
            return False

    except Exception as e:
        logger.warning(f"⚠️ SKIP: API not running (can start with: uvicorn src.app.api.main:app --reload)")
        return None

def test_10_api_sensors_endpoint(self):
    """TEST 10: API /api/v1/sensors/latest endpoint."""
    if self.skip_api:
        logger.info("\n🔊 Skipping API tests (--skip-api flag)")
        return None

    logger.info("\n" + "="*80)
    logger.info("TEST 10: API Sensors Latest Endpoint")
    logger.info("="*80)

    try:
        result = subprocess.run(
            ["curl", "-s", "http://localhost:8000/api/v1/sensors/latest"],

```

```

        capture_output=True,
        text=True,
        timeout=10
    )

    if '"readings"' in result.stdout or "mote" in result.stdout:
        logger.info("✅ PASS: API sensors endpoint responds")
        self.passed_tests.append("API sensors endpoint")
        return True
    else:
        logger.error(f"❌ FAIL: No response from sensors endpoint")
        self.failed_tests.append("API sensors endpoint")
        return False

except Exception as e:
    logger.warning(f"⚠️ SKIP: API not responding")
    return None

# ===== PHASE 5: ML TESTS =====

def test_11_ml_model_exists(self):
    """TEST 11: ML model artifacts exist."""
    logger.info("\n" + "="*80)
    logger.info("TEST 11: ML Model Artifacts")
    logger.info("="*80)

    try:
        model_dir = Path(__file__).parent.parent / "src/app/ml/models"

        if not model_dir.exists():
            logger.error("❌ FAIL: ML models directory doesn't exist")
            self.failed_tests.append("ML model artifacts")
            return False

        model_files = list(model_dir.glob("anomaly_detector_*.pkl"))

        if model_files:
            logger.info(f"✅ PASS: Found {len(model_files)} model artifacts")

            for model_file in model_files[:3]:
                logger.info(f"    {model_file.name}")
                self.passed_tests.append("ML model artifacts")
            return True
        else:
            logger.warning(f"⚠️ SKIP: No model artifacts found (not trained yet)")
            return None

    except Exception as e:

```

```

        logger.warning(f"⚠️ SKIP: {e}")
        return None

# ===== PHASE 6: DASHBOARD TESTS =====

def test_12_dashboard_app_exists(self):
    """TEST 12: Streamlit dashboard app exists."""
    logger.info("\n" + "="*80)
    logger.info("TEST 12: Streamlit Dashboard App")
    logger.info("="*80)

    try:
        dashboard_file = Path(__file__).parent.parent /
"src/app/dashboard/app.py"

        if dashboard_file.exists():
            logger.info("✅ PASS: Dashboard app.py exists")
            self.passed_tests.append("Dashboard app")
            return True
        else:
            logger.warning("⚠️ SKIP: Dashboard app not created yet")
            return None

    except Exception as e:
        logger.warning(f"⚠️ SKIP: {e}")
        return None

# ===== MAIN TEST RUNNER =====

def run_all_tests(self):
    """Run all tests."""
    logger.info("="*80)
    logger.info("SIEIS FULL PIPELINE TEST SUITE")
    logger.info("="*80)
    logger.info(f"Time: {datetime.now()}")
    logger.info(f"Verbose: {self.verbose}")
    logger.info(f"Skip API: {self.skip_api}")
    logger.info("="*80)

    # Phase 1: Docker
    self.test_1_docker_containers_running()
    self.test_2_influxdb_connectivity()
    self.test_3_minio_connectivity()

    # Phase 2: Data Quality
    self.test_4_influxdb_has_incremental_data()
    self.test_5_mote_count_correct()
    self.test_6_sensor_values_valid()
    self.test_7_minio_has_data()

```

```

# Phase 3: Recent Data
self.test_8_recent_data_within_1_hour()

# Phase 4: API
self.test_9_api_health_endpoint()
self.test_10_api_sensors_endpoint()

# Phase 5: ML
self.test_11_ml_model_exists()

# Phase 6: Dashboard
self.test_12_dashboard_app_exists()

# Print summary
self.print_summary()

def print_summary(self):
    """Print test summary."""
    logger.info("\n" + "="*80)
    logger.info("TEST SUMMARY")
    logger.info("="*80)

    total = len(self.passed_tests) + len(self.failed_tests)

    logger.info(f"\n✅ PASSED: {len(self.passed_tests)}")
    for test in self.passed_tests:
        logger.info(f"    ✅ {test}")

    if self.failed_tests:
        logger.info(f"\n❌ FAILED: {len(self.failed_tests)}")
        for test in self.failed_tests:
            logger.error(f"    ❌ {test}")

    logger.info("\n" + "-"*80)

    if len(self.failed_tests) == 0:
        logger.info(f"🎉 ALL TESTS PASSED!")
        logger.info(f"({len(self.passed_tests)}/{total})")
        return True
    else:
        logger.error(f"⚠️ {len(self.failed_tests)} test(s) failed")
        logger.info(f"({len(self.passed_tests)}/{total}) passed")
        return False

def main():
    """Main entry point."""
    parser = argparse.ArgumentParser(

```

```

        description='Run full SIEIS pipeline test suite'
    )
    parser.add_argument(
        '--verbose',
        action='store_true',
        help='Verbose output'
    )
    parser.add_argument(
        '--skip-api',
        action='store_true',
        help='Skip API tests (if API not running)'
    )

    args = parser.parse_args()

    tester = FullPipelineTest(verbose=args.verbose, skip_api=args.skip_api)
    success = tester.run_all_tests()

    sys.exit(0 if success else 1)

if __name__ == "__main__":
    main()

```

tests/test_full_pipeline.py

"""Complete end-to-end pipeline test: Simulator → Kafka → Consumer → InfluxDB + MinIO."""

```

import sys
from pathlib import Path
sys.path.insert(0, str(Path(__file__).parent.parent))

import time
import json
import subprocess
import io
from datetime import datetime
import pandas as pd
from kafka import KafkaConsumer
from influxdb_client import InfluxDBClient
from minio import Minio
from tests.test_config import (
    KAFKA_BROKER,
    KAFKA_TOPIC,
    INFLUX_URL,
    INFLUX_TOKEN,
    INFLUX_ORG,
    INFLUX_BUCKET,

```

```

MINIO_ENDPOINT,
MINIO_ACCESS_KEY,
MINIO_SECRET_KEY,
MINIO_BUCKET,
MINIO_SECURE
)

def test_source_data_count():
    """Test 0: Count records in source data file for today."""
    print("\n" + "="*80)
    print("TEST 0: Source Data Count for Today")
    print("="*80)

    try:
        data_file = Path(__file__).parent.parent / "data" / "processed" /
"realtime_data.txt"

        if not data_file.exists():
            print(f"⚠️ Source data file not found: {data_file}")
            return False

        print(f"📖 Reading source data: {data_file.name}")

        # Read the 9-column space-delimited file
        df = pd.read_csv(
            data_file,
            sep=r'\s+',
            header=None,
            names=['date', 'time', 'epoch', 'moteid', 'temperature',
                  'humidity', 'light', 'voltage', 'updated_timestamp']
        )

        print(f"✅ Total records in file: {len(df):,}")

        # Parse updated_timestamp and filter to today
        # Use format='ISO8601' to handle timestamps with and without
microseconds
        df['updated_timestamp'] = pd.to_datetime(df['updated_timestamp'],
format='ISO8601')
        today = datetime.now().date()
        today_df = df[df['updated_timestamp'].dt.date == today].copy()

        unique_motes = today_df['moteid'].nunique()

        print(f"✅ Records for today ({today}): {len(today_df):,}")
        print(f"✅ Unique motes for today: {unique_motes}")

        if len(today_df) > 0:

```

```

        # Show mote distribution
        mote_counts = today_df['moteid'].value_counts().head(5)
        print(f"\n📊 Top 5 motes by record count:")
        for mote_id, count in mote_counts.items():
            print(f"    Mote {mote_id}: {count:,} records")

        # Show time range
        min_time = today_df['updated_timestamp'].min()
        max_time = today_df['updated_timestamp'].max()
        print(f"\n🕒 Time range: {min_time} to {max_time}")

        print(f"\n✅ Source data available for testing")
        return True
    else:
        print(f"⚠️ No records found for today in source data")
        print(f"    This may indicate the data needs to be regenerated")
        print(f"    Run: python")
data/realtime_mapping/transform_to_realtime.py")
        return False

except Exception as e:
    print(f"❌ Failed to read source data: {e}")
    import traceback
    traceback.print_exc()
    return False

def test_docker_containers_running():
    """Test 1: Verify all Docker containers are running."""
    print("\n" + "="*80)
    print("TEST 1: Docker Containers Running")
    print("="*80)

    required_containers = [
        "sieis-redpanda",
        "sieis-influxdb3",
        "sieis-minio",
        "sieis-simulator",
        "sieis-consumer"
    ]

    try:
        result = subprocess.run(
            ["docker", "ps", "--format", "{{.Names}}"],
            capture_output=True,
            text=True,
            check=True
        )

```

```

running_containers = result.stdout.strip().split('\n')

all_running = True
for container in required_containers:
    if container in running_containers:
        print(f"   ✅ {container}")
    else:
        print(f"   ❌ {container} NOT RUNNING")
        all_running = False

    return all_running
except Exception as e:
    print(f"❌ Failed to check containers: {e}")
    return False

def test_simulator_producing():
    """Test 2: Verify simulator is producing messages to Kafka."""
    print("\n" + "="*80)
    print("TEST 2: Simulator Producing Messages")
    print("="*80)

    try:
        consumer = KafkaConsumer(
            KAFKA_TOPIC,
            bootstrap_servers=[KAFKA_BROKER],
            auto_offset_reset='earliest',
            consumer_timeout_ms=10000,
            value_deserializer=lambda m: json.loads(m.decode('utf-8'))
        )

        print(f"🔍 Listening for messages on {KAFKA_TOPIC}...")

        message_count = 0
        mote_ids = set()
        sample_message = None

        for message in consumer:
            message_count += 1
            msg_data = message.value
            mote_ids.add(msg_data.get('mote_id'))

            if message_count == 1:
                sample_message = msg_data

            if message_count >= 10:
                break

        consumer.close()

```



```

        if sample_message:
            print(f"\n📬 First message received:")
            print(f"    Mote ID: {sample_message.get('mote_id')}")
            print(f"    Temperature: {sample_message.get('temperature')}")
            print(f"    Updated Timestamp: {sample_message.get('updated_timestamp')}")

        print(f"\n✅ Received {message_count} messages from {len(mote_ids)}
motes")

        if message_count > 0:
            print("✅ Simulator is producing messages")
            return True
        else:
            print("❌ No messages received from simulator")
            return False
    except Exception as e:
        print(f"❌ Failed: {e}")
        return False

def test_consumer_writing_influxdb():
    """Test 3: Verify consumer is writing to InfluxDB."""
    print("\n" + "="*80)
    print("TEST 3: Consumer Writing to InfluxDB")
    print("="*80)

    try:
        client = InfluxDBClient(url=INFLUX_URL, token=INFLUX_TOKEN,
org=INFLUX_ORG, timeout=10000)
        query_api = client.query_api()

        query = f'''
from(bucket: "{INFLUX_BUCKET}")
  |> range(start: -24h)
  |> filter(fn: (r) => r._measurement == "sensor_reading")
  |> limit(n: 10)
'''

        print("🔍 Querying InfluxDB for recent sensor data...")
        tables = query_api.query(query, org=INFLUX_ORG)

        record_count = 0
        mote_ids = set()
        fields = set()

        for table in tables:
            for record in table.records:

```

```

        record_count += 1
        mote_id = record.values.get("mote_id")
        field = record.get_field()

        if mote_id:
            mote_ids.add(mote_id)
        if field:
            fields.add(field)

    print(f"✅ Found {record_count} records")
    print(f"✅ Motes: {sorted(list(mote_ids)[:10])}...")
    print(f"✅ Fields: {fields}")

    client.close()

    if record_count > 0:
        return True
    else:
        print("⚠️ No data in InfluxDB yet (may need more time)")
        return False
except Exception as e:
    print(f"❌ Failed: {e}")
    return False


def test_consumer_writing_minio():
    """Test 4: Verify consumer is writing Parquet files to MinIO."""
    print("\n" + "="*80)
    print("TEST 4: Consumer Writing to MinIO")
    print("="*80)

    try:
        client = Minio(
            endpoint=MINIO_ENDPOINT,
            access_key=MINIO_ACCESS_KEY,
            secret_key=MINIO_SECRET_KEY,
            secure=MINIO_SECURE
        )

        print(f"🔍 Checking MinIO bucket: {MINIO_BUCKET}")

        # Limit object listing to prevent timeout with thousands of files
        objects = []
        for obj in client.list_objects(MINIO_BUCKET, prefix="year=",
recursive=True):
            objects.append(obj)
            if len(objects) >= 100: # Check first 100 files only
                break

```

```

    parquet_files = [obj for obj in objects if
obj.object_name.endswith('.parquet')]

    print(f"✅ Found {len(parquet_files)} Parquet files (checked first 100
objects)")

    if parquet_files:
        print("\n📁 Sample files:")
        for obj in parquet_files[:5]:
            print(f"    {obj.object_name} ({obj.size/1024:.1f} KB)")

        # Skip PyArrow validation - too slow for CI/CD
        # Files are validated by structure and size instead
        print(f"\n✅ Parquet files validated by structure and naming
convention")

        return True
    else:
        print("⚠️ No Parquet files found yet (may need more time)")
        return False
except Exception as e:
    print(f"❌ Failed: {e}")
    return False

def test_data_consistency():
    """Test 5: Verify data consistency between InfluxDB and MinIO."""
    print("\n" + "="*80)
    print("TEST 5: Data Consistency Check")
    print("="*80)

    try:
        influx_client = InfluxDBClient(url=INFLUX_URL, token=INFLUX_TOKEN,
org=INFLUX_ORG, timeout=10000)
        query_api = influx_client.query_api()

        query = f'''
from(bucket: "{INFLUX_BUCKET}")
  |> range(start: -24h)
  |> filter(fn: (r) => r._measurement == "sensor_reading")
  |> keep(columns: ["mote_id"])
  |> distinct(column: "mote_id")
'''

        print("🔍 Querying InfluxDB for mote IDs...")
        tables = query_api.query(query, org=INFLUX_ORG)
        influx_motes = set()
        for table in tables:
            for record in table.records:

```

```

        influx_motes.add(record.values.get("mote_id"))

    print(f"✅ InfluxDB has data from {len(influx_motes)} motes")

    minio_client = Minio(
        endpoint=MINIO_ENDPOINT,
        access_key=MINIO_ACCESS_KEY,
        secret_key=MINIO_SECRET_KEY,
        secure=MINIO_SECURE
    )

    # Limit object listing to prevent timeout with thousands of files
    objects = []
    for obj in minio_client.list_objects(MINIO_BUCKET, prefix="year=",
recursive=True):
        objects.append(obj)
        if len(objects) >= 500: # Check first 500 files for mote IDs
            break

    minio_motes = set()
    for obj in objects:
        if "mote_id=" in obj.object_name:
            parts = obj.object_name.split("mote_id=")
            if len(parts) > 1:
                mote_id = parts[1].split("/")[0]
                minio_motes.add(mote_id)

    print(f"✅ MinIO has data from {len(minio_motes)} motes (checked first
{len(objects)} files)")

    if influx_motes and minio_motes:
        common_motes = influx_motes & minio_motes
        print(f"✅ {len(common_motes)} motes present in both systems")

        if len(common_motes) > 0:
            print("✅ Data consistency verified")
            return True
        else:
            print("⚠️ No common motes (may need more time)")
            return False
    else:
        print("⚠️ Insufficient data for comparison")
        return False

    influx_client.close()
except Exception as e:
    print(f"❌ Failed: {e}")
    return False

```

```

def main():
    """Run full end-to-end pipeline test."""
    print("\n" + "🚀" * 40)
    print("FULL END-TO-END PIPELINE TEST")
    print("🚀" * 40)
    print("\nArchitecture:")
    print("  CSV → Simulator Container → Redpanda (Kafka)")
    print("          ↓")
    print("  Consumer Container")
    print("    📡           📡")
    print("  InfluxDB       MinIO")
    print("  (hot path)     (cold path)")

    results = []

    # First check source data
    results.append(("Source Data Count", test_source_data_count()))

    results.append(("Docker Containers", test_docker_containers_running()))

    if results[-1][1]:
        print("\n⌚ Waiting 5 seconds for services to stabilize...")
        time.sleep(5)

        results.append(("Simulator Producing", test_simulator_producing()))

        print("\n⌚ Waiting 20 seconds for consumer to process and write data...")
        time.sleep(20)

        results.append(("Consumer → InfluxDB",
test_consumer_writing_influxdb()))
        results.append(("Consumer → MinIO", test_consumer_writing_minio()))
        results.append(("Data Consistency", test_data_consistency()))

    print("\n" + "="*80)
    print("TEST SUMMARY")
    print("="*80)

    passed = sum(1 for _, result in results if result)
    total = len(results)

    for test_name, result in results:
        status = "✅ PASS" if result else "❌ FAIL"
        print(f"{status}: {test_name}")

```

```

print(f"\n🇩🇪 Total: {passed}/{total} tests passed")

if passed == total:
    print("\n🎉🎉🎉 FULL PIPELINE TEST PASSED! 🎉🎉🎉")
    print("\n✅ All containers are working correctly")
    print("✅ Data is flowing from Simulator → Kafka → Consumer → InfluxDB
+ MinIO")
    return 0
else:
    print(f"\n⚠️ {total - passed} test(s) failed")
    print("\nTroubleshooting:")
    print(" 1. Check container logs: docker-compose logs -f")
    print(" 2. Verify services are healthy: docker ps")
    print(" 3. Check consumer logs: docker-compose logs consumer")
    return 1

if __name__ == "__main__":
    sys.exit(main())

```